

# Probabilistic Modeling for Novelty Detection with Applications to Fraud Identification

---

**Rémi Domingues**

**PhD Defense**

**EURECOM - Sorbonne University**

January 29<sup>th</sup> 2019

Advisor: Maurizio Filippone

Co-advisor: Pietro Michiardi

# Motivations - The Amadeus use case



## Motivations - The Amadeus use case



# Motivations - Fraud detection



Compromised user  
accounts



Fraudulent  
bookings



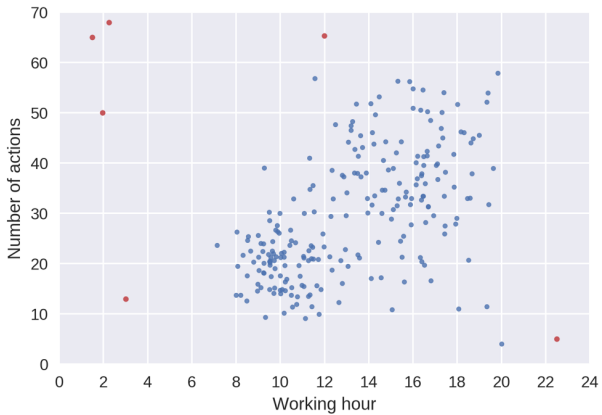
Payment frauds



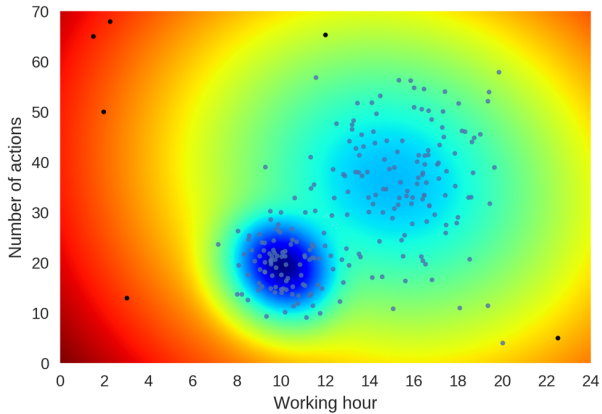
Malicious bots

- Supervised learning – The class imbalance problem
- Unsupervised learning – 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

# Novelty detection



# Novelty detection



# Industrial constraints and research challenges

Proactive, unlabelled data	Novelty detection
Anomalies in training data	Robustness
Continuous scoring	Probabilistic method
Numerical & categorical data	Variational learning of joint distributions
Scalable and distributed	Mini-batch learning
White-box model	Interpretable
Little tuning	Nonparametric



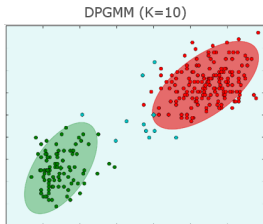


# Dirichlet Process Mixture Model

---

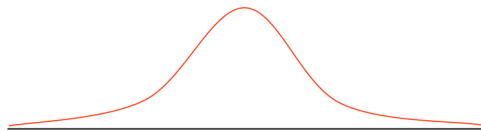
# Dirichlet Process Mixture Model (DPMM)

- **Weighted mixture of multivariate distributions in the exponential family**
- **Nonparametric Bayesian method:**  
infinite-dimensional parameter space
- Dirichlet Process as nonparametric prior
- **A product of exponential-family distributions is in the exponential-family**
- Probabilistic, mini-batch training, categorical support, clustering

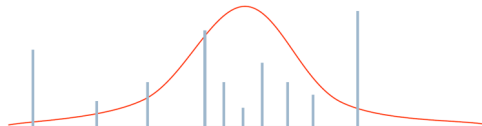


# Dirichlet Process

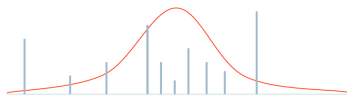
- Bayesian nonparametric model
- **Distribution over distributions**
- Consider a Gaussian  $G_0$ :



- $G \sim DP(\alpha, G_0)$

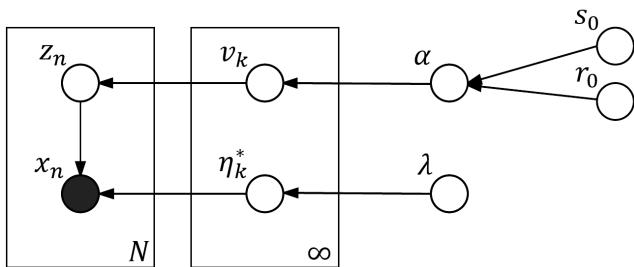


# Stick-breaking process

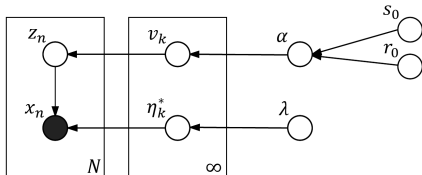


- Constructive way of forming  $G$
- Weights  $\pi_k(\mathbf{v}) = \mathbf{v}_k \prod_{j=1}^{k-1} (1 - v_j)$ , with  $\mathbf{v}_k \sim \text{Beta}(1, \alpha)$
- $G \sim DP(\alpha, G_0) \Leftrightarrow G = \prod_{k=1}^{\infty} \pi_k \delta_{\theta_k}$ 
  - $\theta_k^* \sim G_0$
  - $\delta_{\theta_k}$  is the indicator function which evaluates to zero everywhere, except for  $\delta_{\theta_k}(\theta_k) = 1$

# Dirichlet Process Mixture Model



# Dirichlet Process Mixture Model



1. Draw  $\alpha | s_0, r_0 \sim \Gamma(s_0, r_0)$
2. Draw the **stick length**  $\mathbf{v}_k | \alpha \sim \text{Beta}(1, \alpha)$ ,  
yielding the **mixing weights**  $\pi_k(\mathbf{v}) = \mathbf{v}_k \prod_{j=1}^{k-1} (1 - v_j)$
3. Draw **component**  $\eta_k^* | \lambda \sim G_0$ ,  
with  $G_0$  conjugate prior in the exponential family, e.g.  $p(\mathbf{X} | \eta^*)$   
multivariate normal,  $G_0$  Normal-wishart
4. **Assign** the data to the components:  $z_n | \mathbf{v} \sim \text{Mult}(\pi(\mathbf{v}))$   
Generate the observations:  $\mathbf{x}_n | z_n \sim p(\mathbf{x}_n | \eta_{z_n}^*)$

# DP mixture inference

- Predictive density:  $p(\mathbf{x}_{N+1}|\mathbf{X}, \theta) = \int p(\mathbf{x}_{N+1}|\mathbf{W})p(\mathbf{W}|\mathbf{X}, \theta)d\mathbf{W}$
- Intractable posterior over the latent variables  $p(\mathbf{W}|\mathbf{X}, \theta)$
- Approximate inference
  - Markov Chain Monte Carlo techniques, e.g. Gibbs sampling
  - Variational Inference
  - *Variational inference is that thing you implement while waiting for your Gibbs sampler to converge.* — David Blei

# Variational inference

- Approximate the **posterior**  $\mathbf{p}$  by a **tractable approximation**  $\mathbf{q}$  with variational parameters

- $q$  is from a family of simpler distributions

$$q(\mathbf{v}, \boldsymbol{\eta}^*, \mathbf{z}, \mathbf{w}) = q_{\alpha, \beta}(\mathbf{v}) \cdot q_{\tau}(\boldsymbol{\eta}^*) \cdot q_r(\mathbf{z}) \cdot q_{g_1, g_2}(\mathbf{w})$$

- $q_{\alpha, \beta}(\mathbf{v})$ : product of **Beta**
  - $q_{\tau}(\boldsymbol{\eta}^*)$ : product of distributions in the **exponential family**
  - $q_r(\mathbf{z})$ : product of **multinomials** on cluster assignment variable  $\mathbf{z}$
  - $q_{g_1, g_2}(\mathbf{w})$ :  $\Gamma$  distribution
- Hyperparameters:  $\boldsymbol{\lambda}$ ,  $\alpha$ ,  $s_0$  and  $r_0$
  - Latent variables:  $\mathbf{v}$ ,  $\boldsymbol{\eta}^*$ ,  $\mathbf{z}$  and  $\mathbf{w}$
  - Variational parameters:  $\alpha_k$ ,  $\beta_k$ ,  $\boldsymbol{\tau}_k$ ,  $r_{nk}$ ,  $g_1$  and  $g_2$



# Variational inference

1. Initialize the model parameters

2. Optimize the **variational parameters** to minimize

$$D_{KL}(q(\mathbf{W})||p(\mathbf{W}|\mathbf{X}, \theta)) = \mathbb{E}_q[\ln q(\mathbf{W})] - \mathbb{E}_q[\ln p(\mathbf{W}, \mathbf{X}|\theta)] + \ln p(\mathbf{X}|\theta)$$

$$\text{Equivalent to maximizing } \ln p(\mathbf{X}|\theta) \geq \mathbb{E}_q[\ln p(\mathbf{W}, \mathbf{X}|\theta)] - \mathbb{E}_q[\ln q(\mathbf{W})]$$

3. Compute the **expectation** of  $p(\mathbf{W}|\mathbf{X}, \theta)$  under  $q(\mathbf{W}|\mathbf{X})$ , e.g.

$$\ln q_{\alpha, \beta}^*(\mathbf{v}) = \mathbb{E}_{\eta^*, \mathbf{z}, \mathbf{w}}[\ln p(\mathbf{X}, \mathbf{v}, \eta^*, \mathbf{z}, \mathbf{w})] + c = \prod_{k=1}^{K-1} \text{Beta}(\alpha_k, \beta_k)$$

4. Compute the **geometric means**

- $\mathbb{E}[\ln \mathbf{v}_k]$ ,  $\mathbb{E}[\ln(1 - \mathbf{v}_k)]$ ,  $\mathbb{E}[\eta^*]$ ,  $\mathbb{E}[-a(\eta^*)]$ ,  $\mathbb{E}[z_{nk}]$ ,  $\mathbb{E}[\mathbf{w}]$  and  $\mathbb{E}[\ln \mathbf{w}]$
- Update the model parameters to maximize the expectation of  $p(\mathbf{W}, \mathbf{X}|\theta)$  under  $q(\mathbf{W}|\mathbf{X})$

- Nondecreasing, used for **convergence monitoring**

$$\begin{aligned}\ln p(\mathbf{X}|\theta) &\geq \mathbb{E}_q[\ln p(\mathbf{W}, \mathbf{X}|\theta)] - \mathbb{E}_q[\ln q(\mathbf{W})] \\ &\geq \mathbb{E}_q[\ln p(\mathbf{X}, \mathbf{z}, \eta^*, \mathbf{v}, \mathbf{w}|\theta)] - \mathbb{E}_q[\ln q(\mathbf{z}, \eta^*, \mathbf{v}, \mathbf{w})] \\ &\geq \mathbb{E}_q[\ln p(\mathbf{X}|\mathbf{z}, \eta^*)] + \mathbb{E}_q[\ln p(\mathbf{z}|\mathbf{v})] + \mathbb{E}_q[\ln p(\eta^*|\lambda)] \\ &\quad + \mathbb{E}_q[\ln p(\mathbf{v}|\mathbf{w})] + \mathbb{E}_q[\ln p(\mathbf{w}|s_0, r_0)] - \mathbb{E}_q[\ln q_{\alpha, \beta}(\mathbf{v})] \\ &\quad - \mathbb{E}_q[\ln q_{\tau}(\eta^*)] - \mathbb{E}_q[\ln q_r(\mathbf{z})] - \mathbb{E}_q[\ln q_{g_1, g_2}(\mathbf{w})]\end{aligned}$$

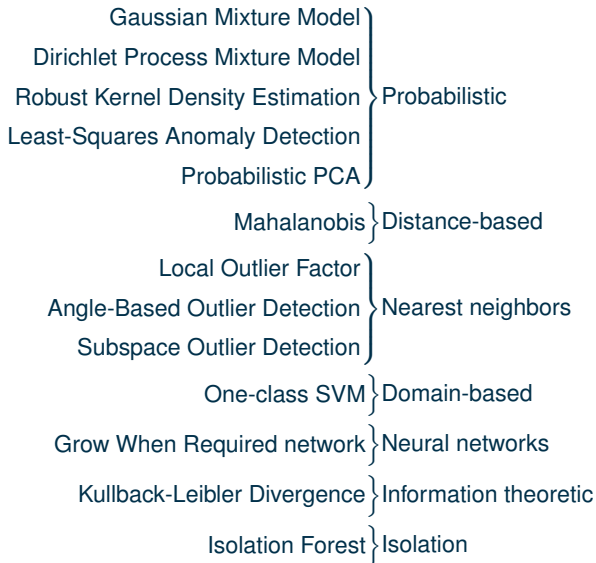
$$\begin{aligned} p(\mathbf{x}_{N+1} | \mathbf{X}, \theta) &= \int \sum_{k=1}^{\infty} \pi_k(\mathbf{v}) p(\mathbf{x}_{N+1} | \eta_k^*) dp(\mathbf{v}, \eta^* | \mathbf{X}, \theta) \\ &\approx \sum_{k=1}^K \mathbb{E}_q[\pi_k(\mathbf{v})] \mathbb{E}_q[p(\mathbf{x}_{N+1} | \eta_k^*)]. \end{aligned} \tag{1}$$

- Analytically, we obtain  $\mathbb{E}_q[\pi_k(\mathbf{v})] = \frac{\alpha_k}{\alpha_k + \beta_k} \prod_{i=1}^{k-1} \left(1 - \frac{\alpha_i}{\alpha_i + \beta_i}\right)$
- **Monte Carlo sampling** is used to estimate the density
  1. Draw  $m$  samples from  $q_{\tau}^*(\eta^*)$
  2. Compute each  $p(\mathbf{x}_{N+1} | \eta^*)$
  3. Average the resulting  $m$  likelihoods

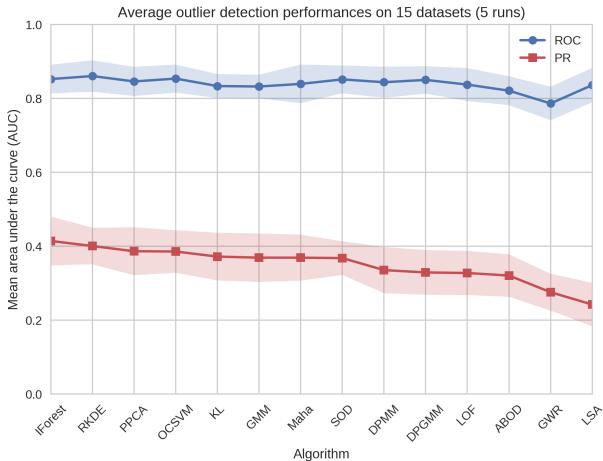
# Experimental survey

---

# Algorithms

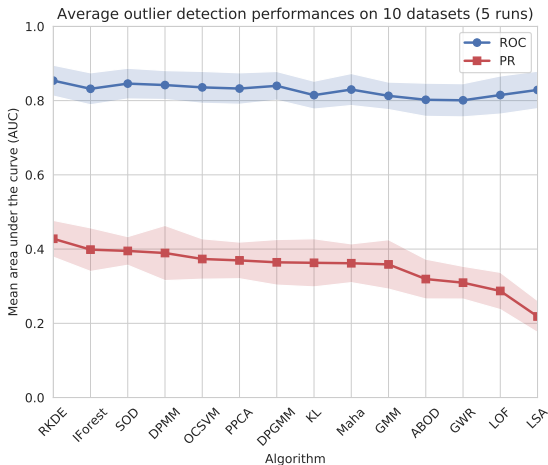


# Results



Algorithm	GMM	DPGMM	DPMM	RKDE	PPCA	LSA	MAHA	LOF	ABOD	SOD	KL	GWR	OCSVM	IFOREST
PR AUC	6	12	9	2	5	14	7	11	10	4	8	13	3	1
ROC AUC	11	4	5	1	7	9	6	10	13	8	12	14	3	2

# Results - No classification datasets



Algorithm	GMM	DPGMM	DPMM	RKDE	PPCA	LSA	MAHA	LOF	ABOD	SOD	KL	GWR	OCSVM	IFOREST
PR AUC	5	4	3	2	6	14	9	12	11	1	8	13	10	7
ROC AUC	13	5	2	1	6	9	7	10	14	4	12	11	3	8

## Area under the ROC and Precision-Recall curves

- 10 frauds, 990 normal transactions (i.e. 1% positives, 99% negatives)
- *Prediction 1*: 5 frauds correctly labelled, all normal transaction correctly labelled
- *Prediction 2*: 5 frauds correctly labelled, 20 normal transactions incorrectly labelled

AUC	ROC	PR
<b>Prediction 1</b>	0.75	0.50
<b>Prediction 2</b>	0.74	0.10

- **The ROC AUC downplays the impact of false positives when negative observations are over-represented**

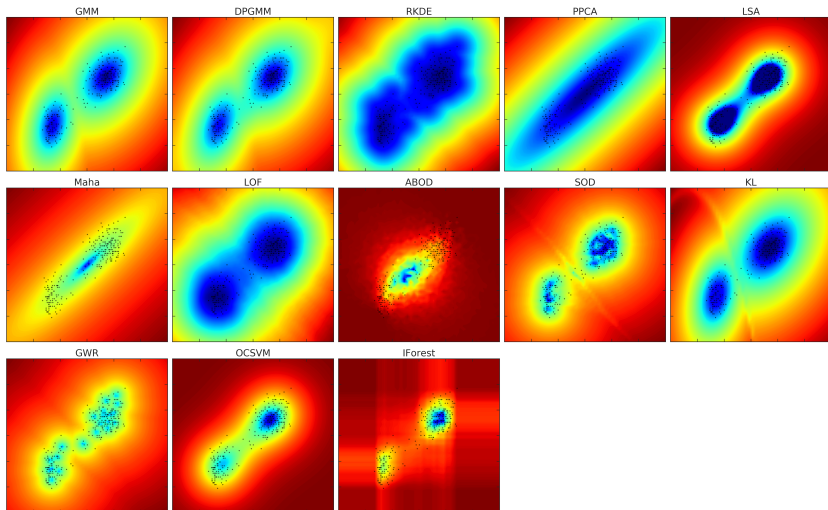


# Scalability

- **Runtime** and **memory** scalability
- Stability, robustness, resistance to the curse of dimensionality
- Datasets of increasing size, dimensionality and noise

Algorithm	Training/prediction time		Mem. usage		Robustness		
	↗ Samples	↗ Features	↗ Samples	↗ Features	↗ Noise	High dim.	Stability
GMM	Low/Low	Medium/Medium	Low	Medium	High	Medium	Medium
BGM	Low/Low	Medium/Medium	Low	Medium	High	Medium	High
DPGMM	Medium/Low	High/High	Low	High	High	High	High
RKDE	High/High	High/High	High	Low	High	High	High
PPCA	Low/Low	High/Low	Low	Low	High	Medium	Medium
LSA	Low/Medium	Low/Low	Medium	Low	Low	Low	Medium
MAHA	Low/Medium	Medium/Low	Low	Medium	Medium	Low	High
LOF	High/High	Low/Low	High	Low	Medium	High	High
ABOD	Low/High	Low/Medium	Low	Low	Medium	Low	Medium
SOD	High/High	Low/Medium	High	Low	Low	High	Medium
KL	Low/Medium	Low/Medium	Low	Medium	High	Medium	High
GWR	Medium/Medium	Medium/Low	Low	Low	Low	High	Medium
OCSVM	High/High	Low/Low	Low	Low	Low	High	High
IFOREST	Low/Medium	Low/Low	Medium	Low	High	High	Medium

# Contours - Old Faithful dataset



# Deep Gaussian Process autoencoder

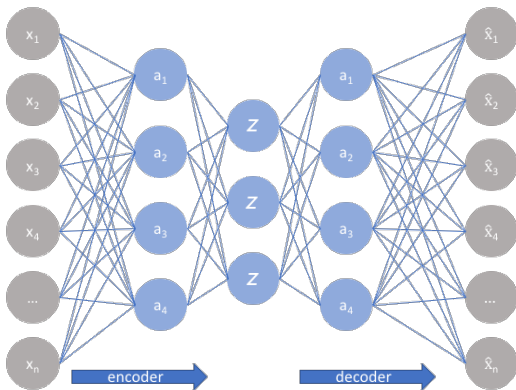
---

# 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

# Autoencoders

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

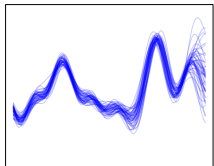


# Deep Gaussian Process autoencoders

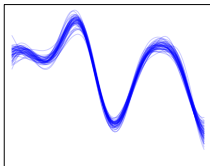
- Deep probabilistic models
- Composition of functions

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

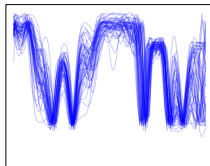
$\mathbf{h}^{(0)}(\mathbf{x})$



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



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



- Inference requires calculating the marginal likelihood:

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

# DGPs with Random Features

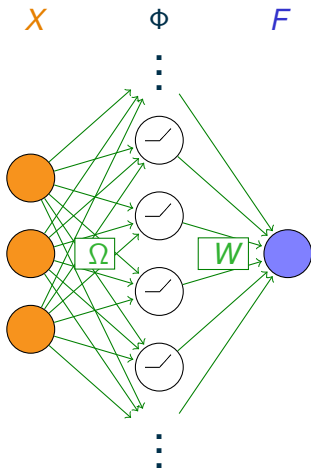
- GPs are single-layered Neural Nets with an infinite number of hidden units

- Weight-space view of a GP

$$F = \Phi W$$

- The priors over the weights are

$$p(W_{.j}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$$





# Random Feature Expansion of Kernels

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

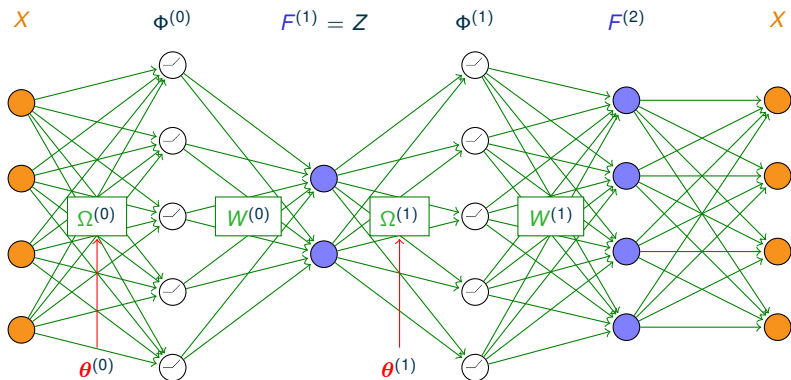
$$\Phi_{\text{RBF}} = \sqrt{\frac{\sigma^2}{N_{\text{RF}}}} [\cos(F\Omega), \sin(F\Omega)] \quad \text{with} \quad p(\Omega_j | \theta) = \mathcal{N}(\mathbf{0}, \Lambda^{-1})$$

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

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

- 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:

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

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

- $D_{\text{KL}}$  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**

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

- **Monte Carlo sampling**

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

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

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

- Reparameterization trick

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

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

- Predictive distribution

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

- Approximation

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

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

# DGP-AE Evaluation

---



- **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)

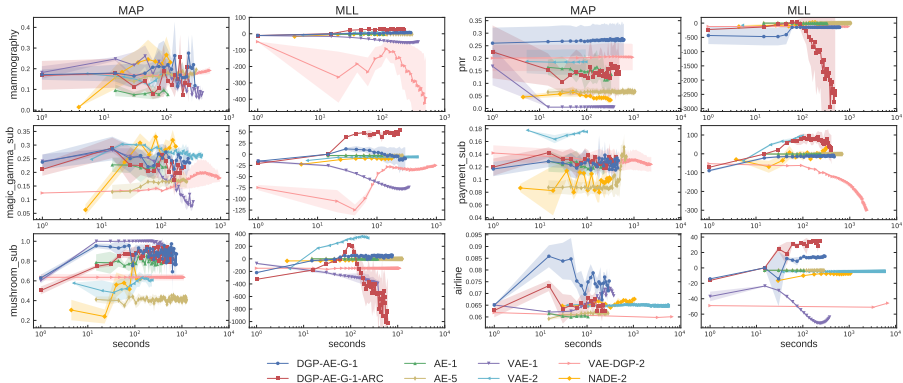
# Method comparison

- 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	<b>0.344</b>	<b>0.338</b>	<b>0.366</b>	<b>0.370</b>	0.284	0.264	0.222	0.262	0.270	0.216	0.336	0.284

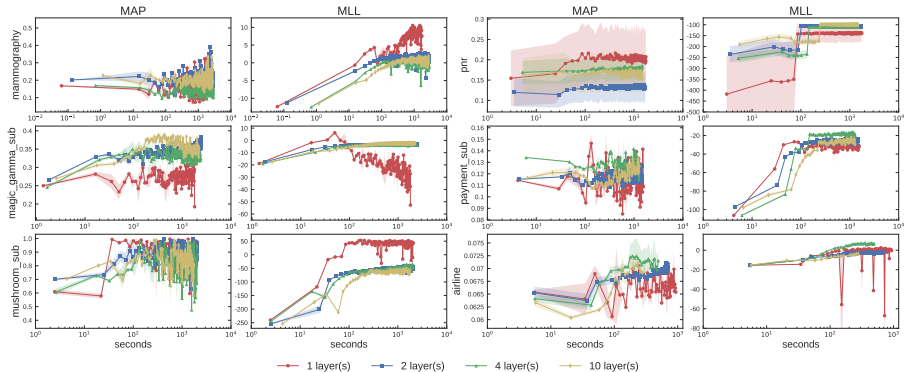
# Convergence monitoring - Networks

- 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

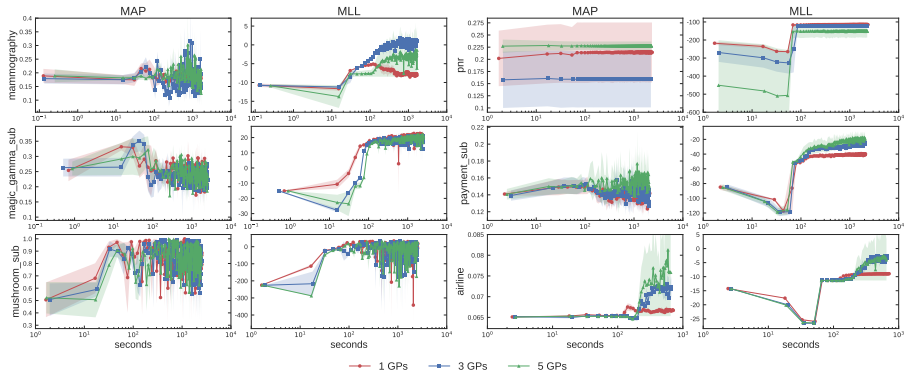
# Convergence monitoring - Depth



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

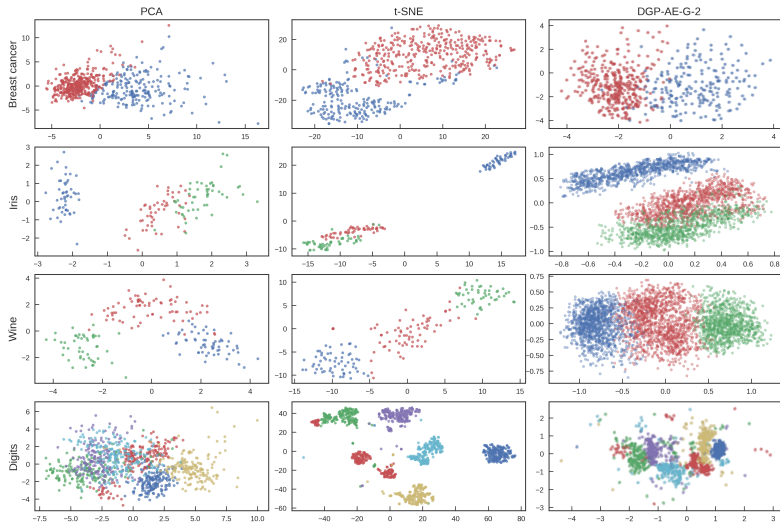
# Convergence monitoring - GPs

- 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

- Novel probabilistic models for novelty detection
  - DPMM
    - Interpretable, fast and accurate modeling of **mixed-type features**
    - Clustering, not suitable for numeric-only data
  - DGP-AE
    - Competitive with SoA and DNN-based novelty detection methods
    - Good **dimensionality reduction** abilities
    - Tractable and scalable inference
    - Suitable to model **mixed-types features**
- Experimental surveys for novelty detection
  - **Numerical, mixed-type** and **temporal** data
  - **No clear winner**
  - Metric comparison
  - Recommendations based on datasets' characteristics
  - Highlighted **scalability** pitfalls



- Generic **benchmarking platform**
- Comparative study used internally
- Thousands of DPMMs running to **raise alerts**
- **Recommendations** for action sequences + **integration** ready

# Research contributions

- **Journals**

R. Domingues, M. Filippone, P. Michiardi, and J. Zouaoui, A comparative evaluation of outlier detection algorithms: experiments and analyses, ***Pattern Recognition***, vol. 74, pp. 406–421, 2018

R. Domingues, P. Michiardi, J. Barlet, and M. Filippone, A comparative evaluation of novelty detection algorithms for discrete sequences, ***Artificial Intelligence Review***, vol. 52, no. 1, 2019

**Under review**

- **Journal & conference**

R. Domingues, P. Michiardi, J. Zouaoui, and M. Filippone, Deep gaussian process autoencoders for novelty detection, ***Machine Learning***, 2018

Presented at ***ECML-PKDD***, 2018

- **Workshop**

R. Domingues, F. Buonora, R. Senesi, and O. Thonnard, An application of unsupervised fraud detection to passenger name records, in *2016 46th Annual IEEE/IFIP International Conference on **Dependable Systems and Networks 35/41 Workshop (DSN-W)***, 2016, pp. 54–59

# Future work

- Mini-batch training for DPMM
- **Generative DGP-AE**
- Model discrete event sequences with **structured DGP-AE**
- Image-based novelty detection
- **Distributed** and **GPU** computing, **streaming** data

**Thank you**

# Exponential family of distributions

- Density
  - $h(\mathbf{x})$  function
  - $\boldsymbol{\eta}^*$  natural parameter
  - $T(\mathbf{x})$  sufficient statistics
  - $a(\boldsymbol{\eta}^*)$  normalization factor

$$p(\mathbf{x}|\boldsymbol{\eta}^*) = h_l(\mathbf{x}) \exp\left(\boldsymbol{\eta}^{*T} T(\mathbf{x}) - a_l(\boldsymbol{\eta}^*)\right) \quad (2)$$

- Conjugate prior, based on the previous likelihood

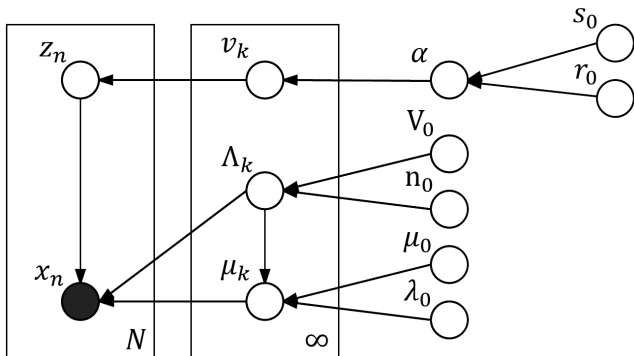
$$p(\boldsymbol{\eta}^*|\boldsymbol{\lambda}) = h_p(\boldsymbol{\eta}^*) \exp\left(\boldsymbol{\lambda}_1^T \boldsymbol{\eta}^* + \lambda_2(-a_l(\boldsymbol{\eta}^*)) - a_p(\boldsymbol{\lambda})\right), \quad (3)$$

- Same dimensionality for  $\boldsymbol{\lambda}$  and  $\boldsymbol{\eta}^*$ ,  $\lambda_2$  is a scalar

- Posterior

$$p(\boldsymbol{\eta}^*|\boldsymbol{\tau}) = h_p(\boldsymbol{\eta}^*) \exp\left(\boldsymbol{\tau}_1^T \boldsymbol{\eta}^* + \tau_2(-a_l(\boldsymbol{\eta}^*)) - a_p(\boldsymbol{\tau})\right). \quad (4)$$

# Dirichlet Process Mixture Model



# Dirichlet Process Mixture Model

- Mean-field variational inference
  - The optimal solution  $q_j^*$  for each of the factors  $q_j$  is:

$$\ln q_j^*(\mathbf{W}_j | \mathbf{X}) = \mathbb{E}_{i \neq j} [\ln p(\mathbf{X}, \mathbf{W})] + \text{const} \quad (5)$$

- Truncated representation of a DP mixture
  - $\pi_k(\mathbf{v}) = \mathbf{v}_k \prod_{j=1}^{k-1} (1 - v_j)$
  - $\pi_k(\mathbf{v}) = 0$  for  $k > K$ , which is achieved by setting  $v_K = 1$
  - $q_{\alpha_K, \beta_K}(v_K = 1) = 1$