Détection d'anomalies non supervisée dans NSL-KDD en utilisant un β -VAE une approche basée sur l'espace latent et l'erreur de reconstruction

Dylan Baptiste*†, Ramla Saddem*, Alexandre Philippot*, François Foyer†

*Université de Reims Champagne-Ardenne, CRESTIC, France †Seckiot, Paris, France







Introduction

- Convergence IT/OT croissante
- Besoin accru en systèmes « Intrusion Detection System »
- Focus sur une approche non supervisée



Objectifs de l'étude

- Détection d'anomalies non supervisé
- Comparaison de deux approches en exploitant un β -VAE :
 - Erreur de reconstruction
 - Distance dans l'espace latent
- Application sur le jeu de données NSL-KDD*

Le modèle β -VAE

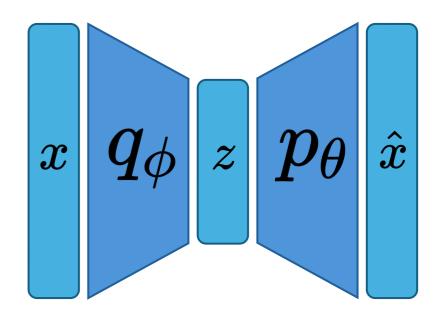
Encodeur

Qui projette une donnée x

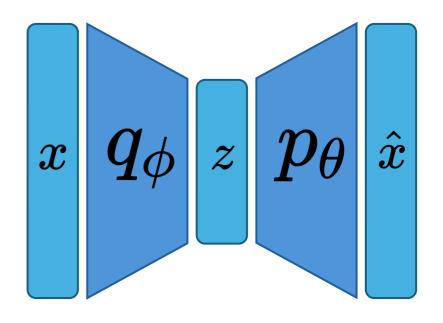
dans un espace latent z

Décodeur p_{θ}

Qui récré la donnée initiale à partir de z



Fonction de coût



$$-\mathbb{E}q_{\phi}(z|x)[\log p_{\theta}(x|z)] + \mathcal{B}_{KL}(q_{\phi}(z|x)||p(z))$$
 Erreur de reconstruction Divergence de Kullback-Leibler

Fonction de coût

Caractéristiques continues

activation linéaire avec mean squared error

$$\mathcal{L}_{\text{cont}} = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{x}_i)^2$$

Caractéristiques booléennes

activation sigmoid avec binarycross-entropy loss

$$\mathcal{L}_{\text{bool}} = -\frac{1}{n} \sum_{i=1}^{n} \left[x_i \log(\hat{x}_i) + (1 - x_i) \log(1 - \hat{x}_i) \right]$$

Caractéristiques catégorielles

activation softmax avec categorical cross-entropy loss

$$\mathcal{L}_{\text{cat}} = -\frac{1}{n} \sum_{i=1}^{n} \sum_{j=1}^{m} x_{ij} \log(\hat{x}_{ij})$$

$$\mathcal{L}_{rec} = \mathcal{L}_{cat} + \mathcal{L}_{bool} + \mathcal{L}_{cont}$$

Méthode 1 : Erreur de reconstruction

Algorithm 1 \mathcal{L}_{rec} -classification

Require: x a sample to classify, (q_{ϕ}, p_{θ}) : a trained β -VAE, τ : the threshold

Ensure: y: classification label: normal or anomalie

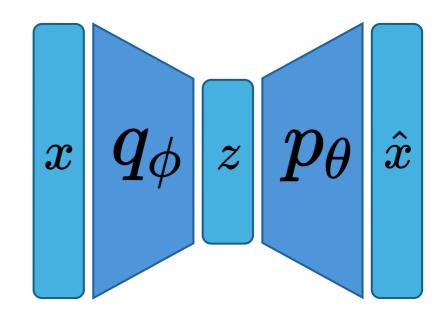
1: $z \sim q_{\phi}(z|x)$

 \triangleright Encoder

2: $\hat{x} \leftarrow p_{\theta}(z)$

▶ Decoder

- 3: if $\mathcal{L}_{\text{rec}}(x,\hat{x}) > \tau$ then
- 4: $y \leftarrow \text{anomalie}$
- 5: **else**
- 6: $y \leftarrow \text{normal}$
- 7: end if
- 8: **return** *y*

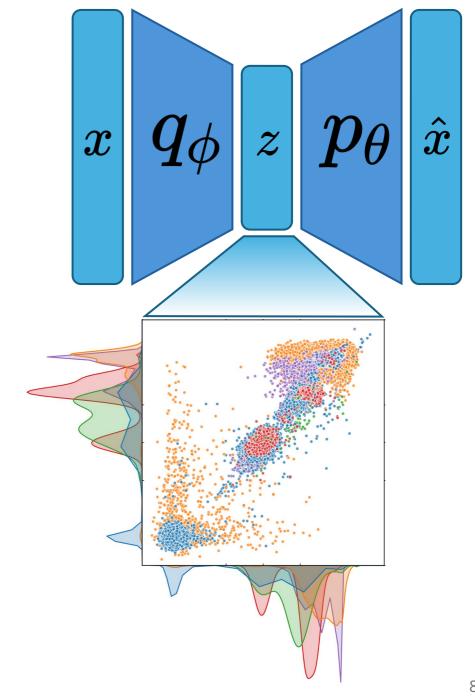


Méthode 2 : Distance dans l'espace latent

On calcule la moyenne des distances de la projection **z** d'une nouvelle observation **x** au **k** plus proches projections d'un ensemble X connu.

$$\mathcal{Z}_k^X(x) = \frac{1}{k} \sum_{j=1}^k \|z - z'_{(j)}\|_2$$

Avec $z \sim q_{\varphi}(x \mid z)$ et $z'_{(j)}$ le j-ème plus proche voisin de z dans l'ensemble des projections de X.



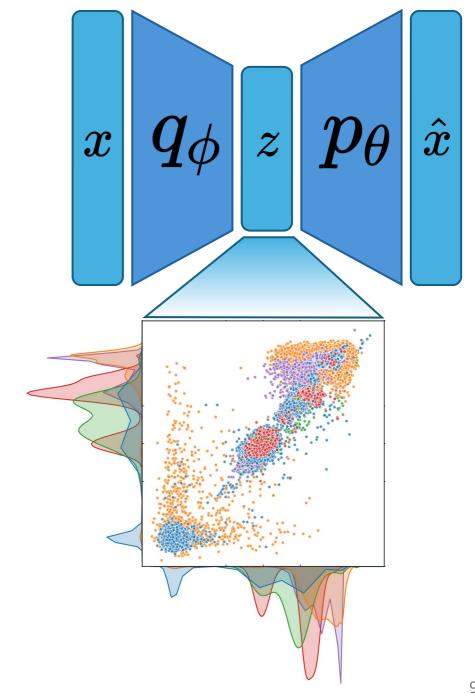
Méthode 2 : Distance dans l'espace latent

Algorithm 2 \mathcal{Z}_k -classification

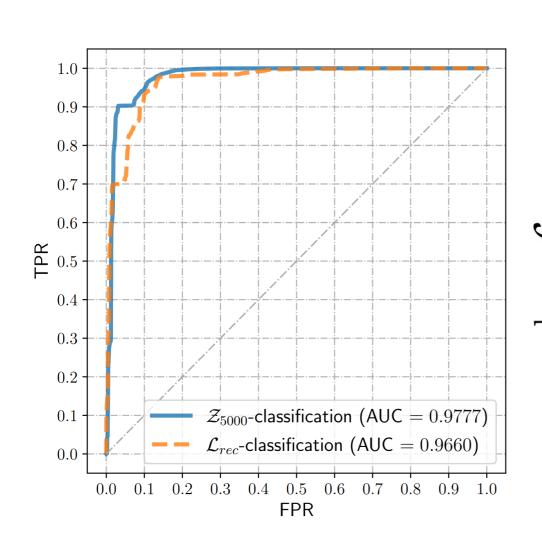
Require: x a sample to classify, (q_{ϕ}, p_{θ}) : a trained β -VAE, X_{train} : the training dataset, k: the number of neighbors, τ : the threshold

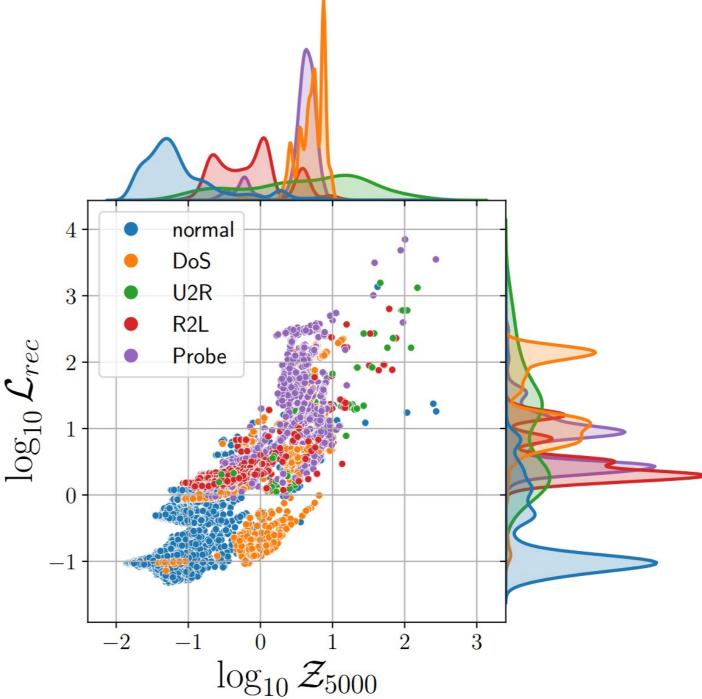
Ensure: y: classification label: normal or anomalie

- 1: $Z_{train} \leftarrow \{z_i \sim q_{\phi}(z|x_i), \forall x_i \in X_{train}\}$
- 2: $z \sim q_{\phi}(z|x)$
- 3: Find the k nearest neighbors $z'_{(1)}, \ldots, z'_{(k)}$ of z in Z_{train}
- 4: if $\mathcal{Z}_k^{X_{train}}(x) > \tau$ then
- 5: $y \leftarrow \text{anomalie}$
- 6: else
- $y \leftarrow \text{normal}$
- 8: end if
- 9: **return** y



Évaluation





Évaluation

Denial of Service (DoS):

Visent à rendre un service ou un système indisponible en le surchargeant de requêtes ou de trafic.

Probe:

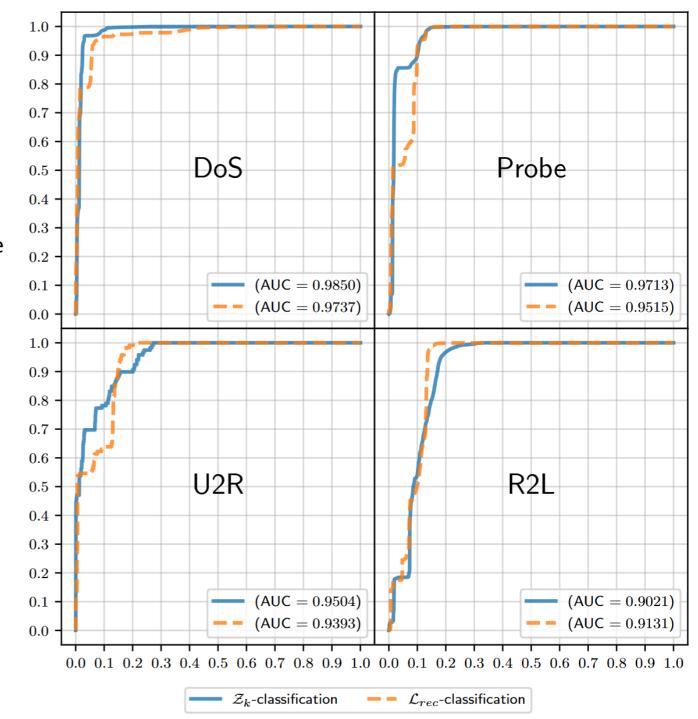
Cherchent à collecter des informations sur le réseau ou les systèmes, comme la cartographie des ports ou des services actifs, en vue d'une attaque ultérieure.

User-to-Root (U2R):

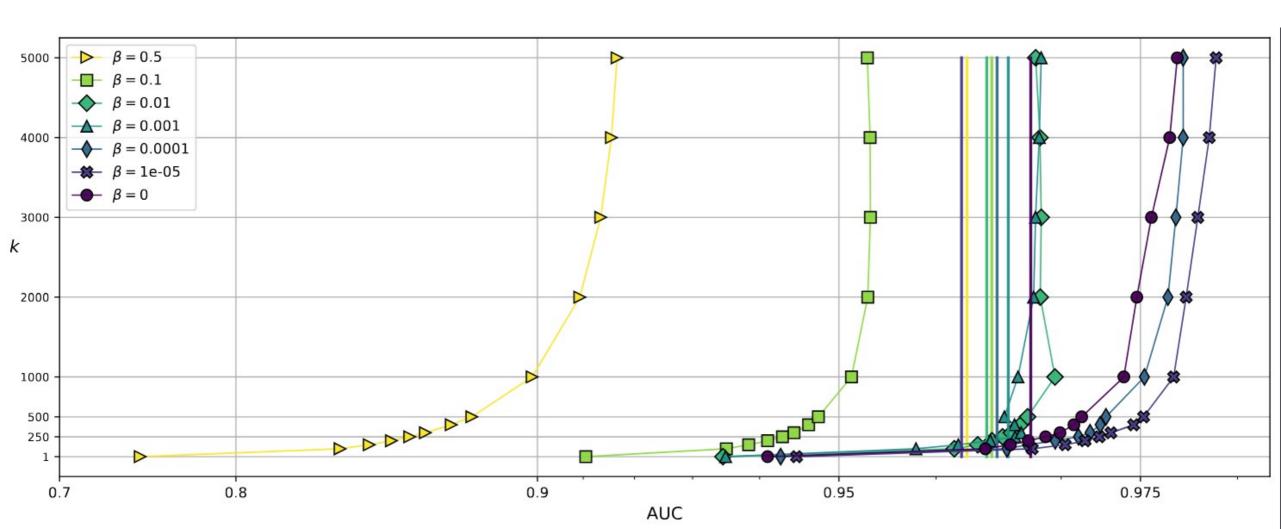
Permettent à un utilisateur malveillant, disposant d'un accès limité à une machine, d'obtenir des privilèges administrateur

Remote-to-Local (R2L):

Permettent à un attaquant distant d'obtenir un accès local non autorisé sur une machine cible



Évaluation



Conclusion

 Les deux méthodes obtiennent d'excellents résultats sur ce jeu de données démontrant ainsi l'efficacité des β-VAE pour la détection non supervisée dans ce contexte

& Perspectives

- Seuils adaptatifs multi-métriques
- Fusion L_{rec} -classification + Z_k -classification
- La Z_k -classification ouvre la voie à une classification multilabel, au-delà d'un simple cadre binaire

