

# GENERATIVE ADVERSARIAL NETWORKS FOR UNSUPERVISED FAULT DETECTION

2018 EUROPEAN CONTROL CONFERENCE (ECC)

by P. Spyridon and Y. S. Boutalis

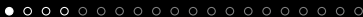
---

Lee Rippon

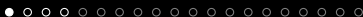
The University of British Columbia  
Chemical and Biological Engineering  
Pulp and Paper Center  
DAIS Lab

October 30, 2019

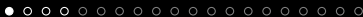
1. Introduction
2. Generative Adversarial Networks
3. Evaluation of GAN Models Based on AEs
4. Hyper-Parameter Selection of GAN Model
5. Unsupervised Fault Detection of Tennessee Eastman Process
6. Conclusions & Future Work



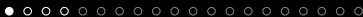
- **Application:** simple fault detection, i.e., identifying deviations from normal operation.



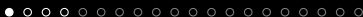
- **Application:** simple fault detection, i.e., identifying deviations from normal operation.
- **Process:** Tennessee Eastman Process. Industrial benchmark data-set for fault detection, diagnosis, isolation and alarm management.



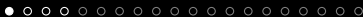
- **Application:** simple fault detection, i.e., identifying deviations from normal operation.
- **Process:** Tennessee Eastman Process. Industrial benchmark data-set for fault detection, diagnosis, isolation and alarm management.
- **Flavor:**



- **Application:** simple fault detection, i.e., identifying deviations from normal operation.
- **Process:** Tennessee Eastman Process. Industrial benchmark data-set for fault detection, diagnosis, isolation and alarm management.
- **Flavor:**
  - Pattern recognition and classification.

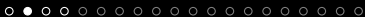


- **Application:** simple fault detection, i.e., identifying deviations from normal operation.
- **Process:** Tennessee Eastman Process. Industrial benchmark data-set for fault detection, diagnosis, isolation and alarm management.
- **Flavor:**
  - Pattern recognition and classification.
  - One-class classification (i.e., anomaly detection).

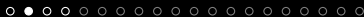


- **Application:** simple fault detection, i.e., identifying deviations from normal operation.
- **Process:** Tennessee Eastman Process. Industrial benchmark data-set for fault detection, diagnosis, isolation and alarm management.
- **Flavor:**
  - Pattern recognition and classification.
  - One-class classification (i.e., anomaly detection).
  - Generative methods.

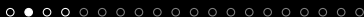




- Fault detection methods are divided into the following three categories:



- Fault detection methods are divided into the following three categories:
  1. **Model based methods** - calculating residuals  $\epsilon = |y - \hat{y}|$ .



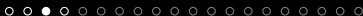
- Fault detection methods are divided into the following three categories:
  1. **Model based methods** - calculating residuals  $\epsilon = |y - \hat{y}|$ .
  2. **Signal processing based methods** - monitoring statistical operations of the system signals.



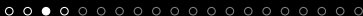
- Fault detection methods are divided into the following three categories:
  1. **Model based methods** - calculating residuals  $\epsilon = |y - \hat{y}|$ .
  2. **Signal processing based methods** - monitoring statistical operations of the system signals.
  3. **Pattern recognition or classification methods** - finding a boundary condition between classes.

- **One-class classification (OCC):** training data includes only one class of data (normal operation). Detect deviations from normal operation.

- **One-class classification (OCC):** training data includes only one class of data (normal operation). Detect deviations from normal operation.
  - **Motivation:** abnormal scenarios are disastrous and rare so the set of faulty data for training is limited.



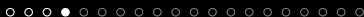
- **One-class classification (OCC):** training data includes only one class of data (normal operation). Detect deviations from normal operation.
  - **Motivation:** abnormal scenarios are disastrous and rare so the set of faulty data for training is limited.
  - **OCC algorithms:** many proposed in literature using neural networks, nearest neighbors, decision trees and bayesian classifiers.



- **One-class classification (OCC):** training data includes only one class of data (normal operation). Detect deviations from normal operation.
  - **Motivation:** abnormal scenarios are disastrous and rare so the set of faulty data for training is limited.
  - **OCC algorithms:** many proposed in literature using neural networks, nearest neighbors, decision trees and bayesian classifiers.
  - **Benchmark:** One-class support vector machines (OCSVM). Popular due to ability to handle non-linear decision boundaries.



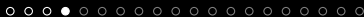
- Typical segue into **deep learning** (DL) methods:



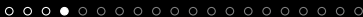
- Typical segue into **deep learning** (DL) methods:
  - Extract features at various levels of representation, achieve a hierarchical representation for learning.

- Typical segue into **deep learning** (DL) methods:
  - Extract features at various levels of representation, achieve a hierarchical representation for learning.
  - Immense success and popularity in computer vision, natural language processing and audio recognition.

- Typical segue into **deep learning** (DL) methods:
  - Extract features at various levels of representation, achieve a hierarchical representation for learning.
  - Immense success and popularity in computer vision, natural language processing and audio recognition.
  - Application in fault and anomaly detection such as acoustic novelty detection using auto-encoders.

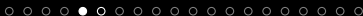


- Typical segue into **deep learning** (DL) methods:
  - Extract features at various levels of representation, achieve a hierarchical representation for learning.
  - Immense success and popularity in computer vision, natural language processing and audio recognition.
  - Application in fault and anomaly detection such as acoustic novelty detection using auto-encoders.
- **Generative models** - estimate the probability distribution function of the training data-set. Popular methods include:



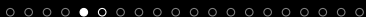
- Typical segue into **deep learning** (DL) methods:
  - Extract features at various levels of representation, achieve a hierarchical representation for learning.
  - Immense success and popularity in computer vision, natural language processing and audio recognition.
  - Application in fault and anomaly detection such as acoustic novelty detection using auto-encoders.
- **Generative models** - estimate the probability distribution function of the training data-set. Popular methods include:
  - Autoregressive models.
  - Variational Auto-Encoders (VAEs).
  - **Generative Adversarial Networks (GANs)**.

# Generative Adversarial Networks



- **Generator** attempts to create realistic samples that fool the discriminator. **Discriminator** is trying to classify if a sample is real or fake. **Non-cooperative game theory.**

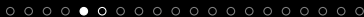
# Generative Adversarial Networks



- **Generator** attempts to create realistic samples that fool the discriminator. **Discriminator** is trying to classify if a sample is real or fake. **Non-cooperative game theory**.
- Two stochastic gradient descent (SGD) algorithms run simultaneously to update parameters of the generator and discriminator.



# Generative Adversarial Networks



- **Generator** attempts to create realistic samples that fool the discriminator. **Discriminator** is trying to classify if a sample is real or fake. **Non-cooperative game theory**.
- Two stochastic gradient descent (SGD) algorithms run simultaneously to update parameters of the generator and discriminator.

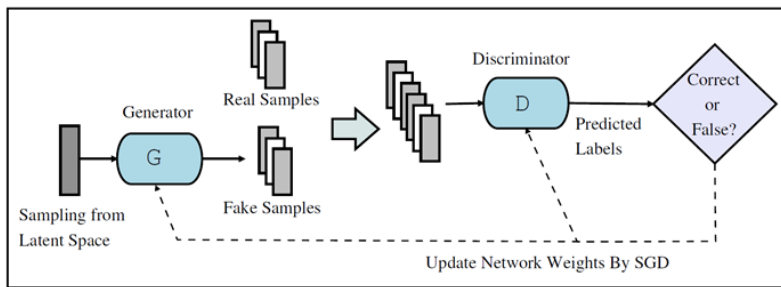


Figure 1: Architecture of Generative Adversarial Network.

- The discriminator output,  $D(x)$ , is the estimated probability that  $x$  comes from the training data rather than the generator.

[illegible]

- $$J^D(\theta^D, \theta^G) = -\frac{1}{2}\mathbb{E}_{x \sim p_{\text{data}}} \log(D(x)) - \frac{1}{2}\mathbb{E}_{z \sim P(z)} \log(1 - D(G(z))). \quad (1)$$

☐ ☐ ☐ ☐ ☒ ☐

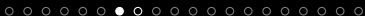
- $$J^D(\theta^D, \theta^G) = -\frac{1}{2}\mathbb{E}_{x \sim p_{\text{data}}} \log(D(x)) - \frac{1}{2}\mathbb{E}_{z \sim P(z)} \log(1 - D(G(z))). \quad (1)$$

- $$J^G(\theta^G, \theta^D) = -\frac{1}{2} \mathbb{E}_{z \sim P(z)} \log(D(G(z))). \quad (2)$$

[illegible]

- Disadvantage of GAN architecture is training instability.

# Generative Adversarial Networks



- Disadvantage of GAN architecture is training instability.
- Solution is not a local minimum of the cost function, instead it is the **Nash equilibrium** of the game, i.e., where each network has no benefit of changing strategy.

# Generative Adversarial Networks



- Disadvantage of GAN architecture is training instability.
- Solution is not a local minimum of the cost function, instead it is the **Nash equilibrium** of the game, i.e., where each network has no benefit of changing strategy.
- **Mode collapse** is a common cause of instability.

# Generative Adversarial Networks

- Disadvantage of GAN architecture is training instability.
- Solution is not a local minimum of the cost function, instead it is the **Nash equilibrium** of the game, i.e., where each network has no benefit of changing strategy.
- **Mode collapse** is a common cause of instability.

Figure 2: John Forbes Nash Jr. played by Russel Crowe in *A Beautiful Mind*, a biographical drama film from 2001).



- Several approaches to enhance GAN stability:

- Several approaches to enhance GAN stability:
  - **Feature matching:** altering generator objective by using intermediate discriminator output.

- Several approaches to enhance GAN stability:
  - **Feature matching:** altering generator objective by using intermediate discriminator output.
  - Addition of **instant noise** to both types of samples before entering discriminator.

- Several approaches to enhance GAN stability:
  - **Feature matching:** altering generator objective by using intermediate discriminator output.
  - Addition of **instant noise** to both types of samples before entering discriminator.
  - Using **dropout** to increase regularization of the generator and avoid over-fitting.

- 10 / 21

- Nearest neighbor evaluation from generated samples.

- **Quantitative evaluation** is an open problem of the GAN framework. Approaches include:
  - Nearest neighbor evaluation from generated samples.
  - Classification performance of the generator.

- Nearest neighbor evaluation from generated samples.
- Classification performance of the generator.
- Comparing two GAN models through a competitive



- **Quantitative evaluation** is an open problem of the GAN framework. Approaches include:
  - Nearest neighbor evaluation from generated samples.
  - Classification performance of the generator.
  - Comparing two GAN models through a competitive game.
  - Evaluating the reconstruction cost of an autoencoder.

- **Autoencoders (AEs):** A feedforward neural network that maps inputs to a hidden latent representation (**encoding**) and then mapping the latent representation back to the original space (**decoding**).

- **Autoencoders (AEs):** A feedforward neural network that maps inputs to a hidden latent representation (**encoding**) and then mapping the latent representation back to the original space (**decoding**).
- Try to minimize the **reconstruction cost** of the input.

- **Autoencoders (AEs):** A feedforward neural network that maps inputs to a hidden latent representation (**encoding**) and then mapping the latent representation back to the original space (**decoding**).
- Try to minimize the **reconstruction cost** of the input.
- **Evaluation approach:** train AE using generator samples and evaluate reconstruction cost  $C^{\text{fake}}$ . Input the original data and evaluate the reconstruction cost  $C^{\text{real}}$ .

- **Autoencoders (AEs):** A feedforward neural network that maps inputs to a hidden latent representation (**encoding**) and then mapping the latent representation back to the original space (**decoding**).
- Try to minimize the **reconstruction cost** of the input.
- **Evaluation approach:** train AE using generator samples and evaluate reconstruction cost  $C^{\text{fake}}$ . Input the original data and evaluate the reconstruction cost  $C^{\text{real}}$ .
- Maximize the evaluation grade  $g^{\text{eval}}$

$$g^{\text{eval}} = \frac{C^{\text{fake}}}{C^{\text{real}}}. \quad (3)$$

- For a set of GAN models, return the model with the highest  $g^{\text{eval}}$ .

- For a set of GAN models, return the model with the highest  $g^{\text{eval}}$ .

---

**Algorithm 1** Calculate the evaluation grade  $g^{\text{eval}}$  of GAN model based on AE

---

1. Train the GAN model using the training dataset of real samples
  2. Generate a batch of fake samples with equal size as the batch of real ones
  3. Train an AE with training dataset the batch of fake samples and estimate  $C^{\text{fake}}$
  4. Calculate the reconstruction cost of AE for the set of real samples  $C^{\text{real}}$ , without training
  5. Return the ratio  $\frac{C^{\text{fake}}}{C^{\text{real}}}$
-

- Generator and discriminator **hyper-parameters** are key considerations to maintain Nash equilibrium.



- Generator and discriminator **hyper-parameters** are key considerations to maintain Nash equilibrium.
- A grid search algorithm is used to select the model hyper-parameters.

- Generator and discriminator **hyper-parameters** are key considerations to maintain Nash equilibrium.
- A grid search algorithm is used to select the model hyper-parameters.
- Final model selection decision is taken by calculating the  $g$ -mean accuracy, i.e.,  $g^{\text{mean}} = \sqrt{a^+ a^-}$  where  $a^+$  and  $a^-$  is the classification accuracy on normal and abnormal data, respectively.

---

**Algorithm 2** Model Selection of GAN for unsupervised classification

---

1. Construct the grid of hyper-parameters
  2. Train  $N$  GAN models, where each model corresponds to a combination of hyper-parameters ( $N$  is the number of hyper-parameter combinations)
  3. For each trained model  $i$ , calculate the evaluation grade  $g_i^{eval}$  using algorithmic procedure 1
  4. Calculate the normalized evaluation grade of each model with the equation  $g_i^{neval} = \frac{g_i^{eval}}{\max_i g_i^{eval}}$
  5. Find the accuracy of the Discriminator on real samples  $a_i^+$  for each model  $i$ .
  6. Find the accuracy of the Discriminator on fake samples  $a_i^-$  for each model  $i$ .
  7. The accuracy on fake samples for each model is updated as  $a_i^- = a_i^- * g_i^{neval}$
  8. The g-mean of each model calculated by the equation  $g_i^{mean} = \sqrt{a_i^+ a_i^-}$
  9. Finally, select the model  $j$  where  $j = \operatorname{argmin}_i g_i^{mean}$
-

- The TE process data-set contains 52 process variables and 21 faulty scenarios.

- The TE process data-set contains 52 process variables and 21 faulty scenarios.
- The normal class training set is 480 samples and the testing set is 960 samples. In this paper these are combined into a normal training set of length 1440.

- The TE process data-set contains 52 process variables and 21 faulty scenarios.
- The normal class training set is 480 samples and the testing set is 960 samples. In this paper these are combined into a normal training set of length 1440.
- Each faulty class has a testing set of 960 samples where the first 160 samples (8 hours) are normal operation.

- Both the generator and discriminator networks have the follow characteristics:

- Both the generator and discriminator networks have the follow characteristics:
  - 5 hidden layers (fixed number of nodes per layer)



- Both the generator and discriminator networks have the follow characteristics:
  - 5 hidden layers (fixed number of nodes per layer)
  - Leaky rectified linear units (ReLU) activation functions.

- Both the generator and discriminator networks have the follow characteristics:
  - 5 hidden layers (fixed number of nodes per layer)
  - Leaky rectified linear units (ReLU) activation functions.
  - Trained with the Adam optimizer.

- Both the generator and discriminator networks have the follow characteristics:
  - 5 hidden layers (fixed number of nodes per layer)
  - Leaky rectified linear units (ReLU) activation functions.
  - Trained with the Adam optimizer.
- The output layer for the generator has linear mappings whereas the output layer of the discriminator uses the softmax function.

- Both the generator and discriminator networks have the follow characteristics:
  - 5 hidden layers (fixed number of nodes per layer)
  - Leaky rectified linear units (ReLU) activation functions.
  - Trained with the Adam optimizer.
- The output layer for the generator has linear mappings whereas the output layer of the discriminator uses the softmax function.
- Additive noise, dropout and feature mapping are used to boost the stability of the generator.

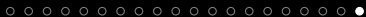
TABLE I  
EVALUATION OF GAN MODELS WITH SEARCH ALGORITHM

Model	Generator Nodes	Discriminator Nodes	Sampling dimension	Evaluating g-mean	Testing g-mean
1	400	600	500	.882	73.76%
2	<b>400</b>	<b>600</b>	<b>1000</b>	<b>.953</b>	<b>75.04%</b>
3	400	800	500	.806	74.66%
4	400	800	1000	.939	75.10%
5	400	1000	500	.835	73.42%
6	400	1000	1000	.922	75.64%
7	500	600	500	.664	69.93%
8	500	600	1000	.567	68.08%
9	500	800	500	.734	73.39%
10	500	800	1000	.623	71.69%
11	500	1000	500	.739	73.39%
12	500	1000	1000	.489	66.04%
13	600	600	500	.687	62.32%
14	600	600	1000	.683	64.98%
15	600	800	500	.695	62.33%
16	600	800	1000	.703	58.40%
17	600	1000	500	.712	62.56%
18	600	1000	1000	.743	67.37%

TABLE II  
FAULT DETECTIONS AND FALSE ALARMS RATES ON TE PROCESS

Fault	Fault Detection Rates		False Alarm Rates		g-means	
	GAN model	OCSVM	GAN model	OCSVM	GAN model	OCSVM
1	<b>99.625%</b>	99.5%	3.125%	<b>0.625%</b>	98.24%	<b>99.44%</b>
2	98.5%	98.5%	1.25%	<b>0%</b>	98.62%	<b>99.25%</b>
3	<b>10.375%</b>	7.625%	<b>7.5%</b>	10.625%	<b>30.98%</b>	26.11%
4	<b>56.25%</b>	50.375%	3.75%	<b>0.625%</b>	<b>73.58%</b>	70.75%
5	<b>32.375%</b>	30.5%	3.75%	<b>0.625%</b>	<b>55.82%</b>	55.05%
6	100%	100%	0%	0%	100%	100%
7	<b>100%</b>	99.625%	1.875%	<b>0%</b>	99.06%	<b>99.81%</b>
8	<b>97.875%</b>	97.375%	1.25%	<b>0%</b>	98.31%	<b>98.68%</b>
9	<b>8.625%</b>	7.125%	<b>9.375%</b>	21.875%	<b>27.96%</b>	23.59%
10	50.875%	<b>53.25%</b>	0.625%	<b>0%</b>	71.10%	<b>72.97%</b>
11	<b>58%</b>	54.75%	2.5%	<b>0.625%</b>	<b>75.20%</b>	73.76%
12	<b>98.75%</b>	98.625%	<b>4.375%</b>	13.125%	<b>97.17%</b>	92.56%
13	<b>95%</b>	94.875%	1.875%	1.875%	<b>96.55%</b>	96.49%
14	100%	100%	1.875%	<b>0%</b>	99.06%	<b>100%</b>
15	12.5%	<b>14%</b>	2.5%	<b>0%</b>	34.91%	<b>37.42%</b>
16	34.375%	<b>36.375%</b>	23.75%	<b>20.625%</b>	51.20%	<b>53.73%</b>
17	<b>91.125%</b>	87.25%	1.875%	<b>0%</b>	<b>94.56%</b>	93.41%
18	<b>90.375%</b>	90.125%	1.875%	<b>0%</b>	94.17%	<b>94.93%</b>
19	<b>11.875%</b>	3.75%	0.625%	0.625%	<b>34.35%</b>	19.30%
20	<b>58.375%</b>	52.75%	0%	0%	<b>76.40%</b>	72.63%
21	<b>49.875%</b>	41.875%	<b>5.625%</b>	6.25%	<b>68.61%</b>	62.66%
Average	<b>64.51%</b>	62.78%	3.77%	<b>3.69%</b>	<b>75.04%</b>	73.45%

- Pri



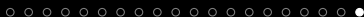
- Conclusions & future work





- Conclusions & future work
- experience

- Conclusions & future work
- experience
- **Models** are anything an agent



- Conclusions & future work
- experience
- **Models** are anything an agent
- **Distributed models**