



Business Intelligence

PUC
RIO

Dijalma Fardin Junior

João R. B. Zoghbi Filho

Rodrigo Pandolfi

*Detecção de Anomalias em Séries Temporais de
Dados de Sistemas de Compressão de Gás*

Monografia de Final de Curso

10/10/2021

***Monografia apresentada ao Departamento de Engenharia Elétrica da
PUC/Rio como parte dos requisitos para a obtenção do título de
Especialização em Business Intelligence.***

Orientadores:

Leonardo Alfredo Forero Mendoza

Dedicatória

Dedico este trabalho aos familiares pelo apoio e compreensão.

Agradecimentos

Aos docentes

RESUMO

Esta monografia apresenta algoritmos estatísticos e de máquinas de aprendizagem no contexto de detecção de anomalias em séries temporais de multivariáveis, aplicados à detecção de falhas de equipamentos, sendo testados em um banco de dados de processos contínuos de um sistema de compressão. Para fomentar a investigação, foi apresentado e discutido um resumo do estado da arte destes métodos de detecção de anomalias. Desta forma, estes métodos foram implementados e avaliados em um programa desenvolvido em Python, utilizando métricas da Matriz de Confusão na avaliação dos mesmos.

ABSTRACT

This work presents statistical and machine learning algorithms in the context of anomaly detection in multivariate time series focused in fault detection of machines, tested in a streaming process dataset of a compressors system. In order to support this investigation, it was presented and discussed, a briefing of the state of the art of these anomaly detection methods. Therefore, these presented methods were implemented and evaluated in a Python Notebook, using Confusion matrix parameters to evaluate them.

Keywords: Outlier Detection; Fault Detection; Machine Learning, Mahalanobis, OC-SVM, Isolation Forest, Autoencoder Neural Network.

SUMÁRIO

1. INTRODUÇÃO	8
1.1 MOTIVAÇÃO	10
1.2 OBJETIVOS DO TRABALHO	10
1.3 DESCRIÇÃO DO TRABALHO	10
2. DESCRIÇÃO DO PROBLEMA	11
2.1 TÉCNICAS DE DETECÇÃO	12
2.2 TRANSFORMAÇÃO DE DADOS	12
2.3 ROTULAÇÃO DE DADOS	12
3.METODOLOGIAS	14
3.1 MÉTODOS DE DETECÇÃO DE ANOMALIA BASEADOS EM ESTATÍSTICA	14
3.1.1 Z SCORE	15
3.1.2 MAHALANOBIS	16
3.2 MÉTODOS DE DETECÇÃO DE ANOMALIA COM APRENDIZAGEM DE MÁQUINA	18
3.2.1 MÁQUINA DE VETORES DE SUPORTE DE CLASSE ÚNICA (ONE CLASS SVM)	18
3.2.2 ISOLATION FOREST	18
3.2.3 REDE NEURAL AUTOENCODER	19
3.4 MÉTRICAS DE DESEMPENHO	21
4. ARQUITETURA DO SISTEMA PROPOSTO	22
5.RESULTADOS	26
6.CONCLUSÕES E TRABALHOS FUTUROS	34
7. REFERÊNCIAS	35

1. INTRODUÇÃO

Em um sistema marítimo de produção de petróleo, do tipo FPSO, Floating, Production, Storage and Offloading, Figura 1, seu sistema de compressão de gás é essencial para assegurar os sistemas produção de petróleo que utilizam: (a) o método de elevação de petróleo, gás lift, que promove o contrafluxo de gás através do anular da coluna de produção, para escoar petróleo da zona de formação à plataforma; e (b) na injeção de gás pressurizado no reservatório para aumentar seu potencial de produção, conforme destacados nos fluxos da Figura 2.

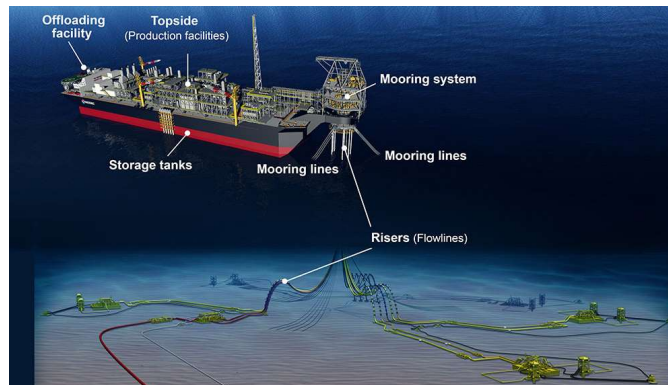


Figura 1: Sistema Flutuante de Produção, Estocagem e Alívio de Petróleo Marítimo e seu arranjo submarino de produção e injeção de fluidos (Fonte: www.modec.com).

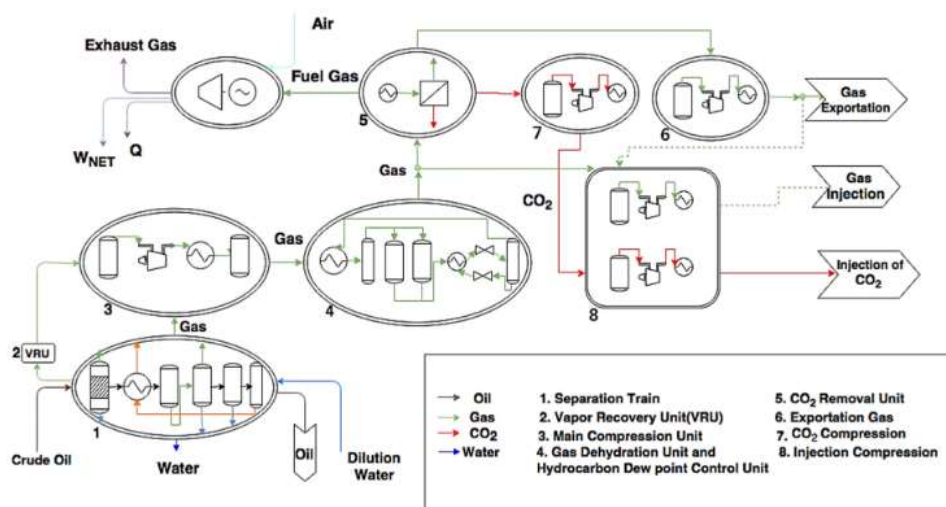


Figura 2: Fluxo de Gás Lift e Gás de Injeção, destacados em verde (Fonte: Bidgoli, 2018).

Além disso, o sistema de compressão fomenta gás para os turbogeradores utilizados na geração elétrica da plataforma, também destacado na Figura 2.

Ou seja, a redução da eficiência ou parada da compressão de gás, pode reduzir ou mesmo interromper a produção dos poços, gerando lucro cessante, e comprometer a geração da planta através de turbo geradores (turbinas à gás ou diesel), embora esta última podendo ser contingenciada pela geração à diesel, porém elevando os custos de produção, devido ao custo do diesel frente ao gás produzido.

Os sistemas de compressão, envolvem um circuito com três estágios principais de compressão que separa o condensado ou líquido residual no gás, com o separador de gás, o compressor que eleva a pressão do gás para os níveis requeridos, e um trocador de calor, que reduz a temperatura

do gás a limites toleráveis para o uso dos consumidores do sistema de compressão, além da válvula antisurge que recircula o gás no compressor, vide Figura 3:

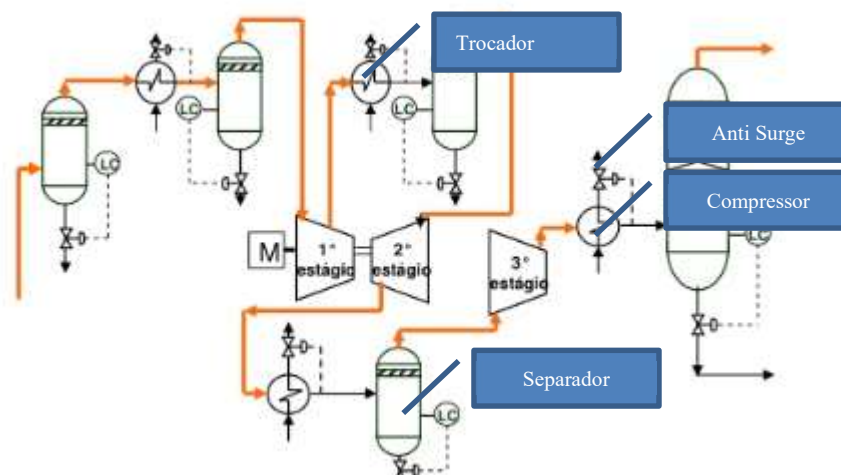


Figura 3: Fluxograma de Sistema de Compressão Principal (Fonte: Carapeto, 2016).

Um sistema de compressão, utiliza o ciclo termodinâmico Brayton para promover trabalho de eixo a sua turbina, conforme a Figura 4, cujo aumento de temperatura e pressão do gás é uma característica inerente ao seu ciclo devido ao ganho de energia durante a compressão combinado com a combustão, sendo arrefecido no trocador de calor.

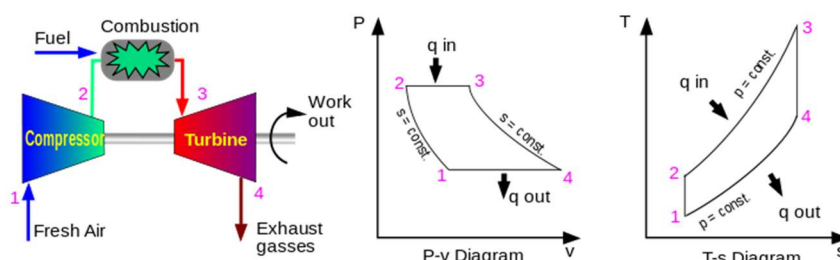


Figura 4: Ciclo de Brayton e diagramas de pressão e temperatura (Fonte: Wikipedia).

Assim, o comportamento variáveis de processo do gás, como pressão, temperatura e vazão refletem o desempenho do sistema de compressão e sua influência direta ou indireta na sua continuidade operacional

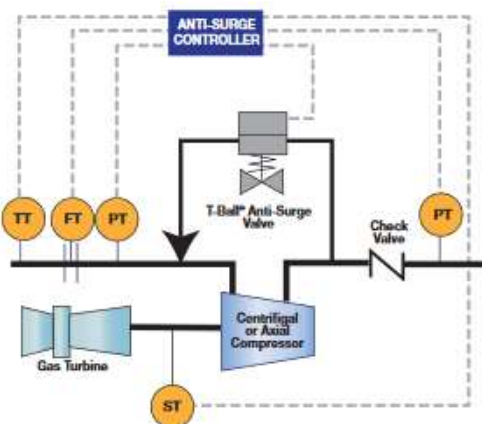


Figura 5: Variáveis de Controle e Intertravamento de Compressores de Gás (Fonte: www.bakerhughes.com)

Dado que, a alteração do comportamento das variáveis de processos do sistema de compressão, indiquem e antecipem o cenário de degradação deste sistema. A seleção e o tratamento destes dados de processo, e a sua modelagem, podem permitir o desenvolvimento de uma ferramenta computacional de detecção de anomalias em séries temporais voltadas para estes sistemas de compressão.

A proposta é então, a partir de um banco de dados históricos de variáveis de processo de temperatura, pressão e vazão de gás em um sistema de compressão, utilizar técnicas de detecção de anomalias, para identificar os momentos em que as mesmas ocorrem, em seguida, testar e validar a metodologia de detecção de anomalias para este sistema de compressão

1.1 MOTIVAÇÃO

A detecção de anomalias quando adotada e implementada para atuar de forma automática e antecipada com base em dados de processo que se comportam como séries temporais, se mostra uma ferramenta de previsão de falhas, ou seja uma evolução da manutenção preditiva vide Kamal e Suganhi (2020) e Janke (2015), antecipando diagnósticos para suporte à decisão das equipes de operação e manutenção para assim mitigar falhas, caracterizando em um processo de manutenção baseado em condição, conforme Figura 6, e assim proporcionando redução de custos e lucro cessante qualquer que seja o ambiente de produção.



Figura 6: Processo de Manutenção baseada em Condição (Brand, 2017).

1.2 OBJETIVOS DO TRABALHO

Desenvolver um algoritmo de detecção de anomalias em dados contínuos de variáveis de processo de um sistema de compressão (séries temporais) utilizando modelos baseados em estatística e aprendizagem de máquinas.

1.3 DESCRIÇÃO DO TRABALHO

Desta forma, o presente trabalho apresenta uma investigação acerca da detecção de anomalias, destacando métodos estatísticos e aprendizagem de máquinas de detecção de anomalias aplicadas em séries temporais, que foram avaliados através de um programa desenvolvido em Python, que testa estes modelos utilizando uma base de dados de processos de um sistema de compressão, considerando a métricas de avaliação Precisão, baseada em Matriz de Confusão.

Assim, nos capítulos seguintes serão discutidos:

- **Capítulo 2, Descrição do Problema**, o qual é discutido as características do problema detecção de anomalias;
- **Capítulo 3, Metodologias**, onde apresenta-se os métodos de detecção de anomalias adotados no desenvolvimento do modelo elaborado deste trabalho;
- **Capítulo 4, Arquitetura**, onde se apresenta o algoritmo da solução proposta para: a análise e tratamento de dados, assim como implementação, otimização de hiperparâmetros através de treino e teste, avaliação dos métodos de detecção de anomalia;
- **Capítulo 5, Testes**, onde apresenta-se os resultados avaliação dos métodos de detecção de anomalia;
- **Capítulo 6, Conclusão**, onde destaca-se os principais aspectos e resultados do processo de revisão bibliográfica, análise e tratamento de dados, desenvolvimento e implementação dos modelos, seus testes e avaliação;

2 DESCRIÇÃO DO PROBLEMA

Conforme citado em Chepoldi (2010), que investigou o tratamento de anomalias de séries temporais, as anomalias são não conformidades identificadas quando comparados ao comportamento normal, vide Figura 7, dentro de um contexto e domínio predeterminado, tais como padrões modificados de dados contínuos, ou descontinuidades observadas na validação de predições, entre outros.

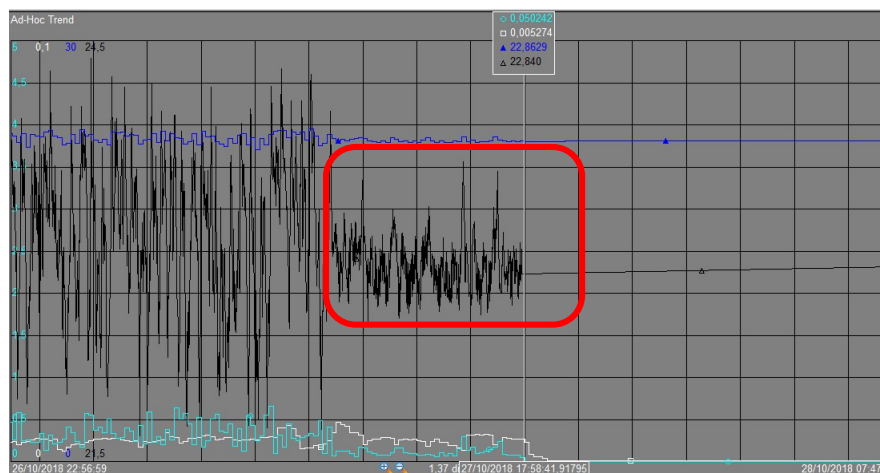


Figura 7: Exemplo de comportamento anômalo destacado em vermelho seguido de parada do equipamento (Fonte: PI Processbook)

Segundo Cheboldi (2010), as técnicas de detecção de anomalias de séries temporais, podem ser classificadas conforme processo de identificação da anomalia ou conforme processo de transformação dos dados (Transformação Dimensional), conforme apresentados no trabalho de Cheboldi, resumido a seguir:

2.1 TÉCNICAS DE DETECÇÃO:

- a) **Baseadas em Janelas Móveis**, onde intervalos de tempo, ou janelas de tempo (subsequências), são definidos em quantidade e extensão para localizar anomalias em uma ou mais janelas;
- b) **Baseadas em Proximidade ou Distância**, que compara a distância entre valores de treino e teste para graduar uma nota de anomalia;
- c) **Baseadas em Predição**, que utiliza um modelo de predição da série temporal (treino), que defasados no tempo são comparados com dados atuais (teste), cuja diferença identifica as anomalias;
- d) **Baseadas em Aprendizagem de Máquinas**, que gera uma série dita escondida, formada pelos dados anômalos mapeados por pesos;
- e) **Baseadas em Segmentação**, no qual a série é fragmentada em séries submetidas à técnica FSA (Fourier Spectrum Analysis), e a probabilidade de transição entre os segmentos são utilizadas para predizer se a natureza é anômala.

2.2 TRANSFORMAÇÃO DE DADOS

Devido a diversidade de alterações associadas às séries, tais como multidimensionalidades, ruídos, escalonamento, entre outros a transformação de dados se faz necessária, e durante este processo as anomalias podem ser detectadas e tratadas em outra dimensão, ou espaço. A seguir são descritos de forma sucinta os processos de Transformação de Séries Temporais:

- a) **Agregação**, que reduzem as dimensões das séries ao comprimi-las;
- b) **Discretização**, que convertem séries temporais em sequencias delimitadas por uma faixa de valores, representadas por letras do alfabeto;
- c) **Baseada no Processamento do Sinal**, conforme o tipo de análise do domínio, seja no tempo, ou frequência, ou ambos, utilizando conforme espectro ou resolução;

2.3 ROTULAÇÃO DE DADOS

A rotulação de dados classifica as instâncias em Normal ou Anômala. O fato de o conjunto de dados ser rotulado, ou não, indicará o procedimento de rotulação a ser utilizado. Dessa maneira as técnicas de detecção de anomalia duvidem-se nas seguintes categorias:

- a). Detecção de anomalia não supervisionada: não necessita de dados de treinamento e normalmente são mais aplicáveis;
- b) Detecção de anomalia semi-supervisionada: as classes normais já estão ao menos definidas. Não há exigência de anomalias rotuladas;

c) Detecção de anomalia supervisionada: há disponibilidade de conjuntos de dados rotulados para as classes normal e anômala.

Assim, as saídas (output) são classificadas de acordo com a técnica empregada: pontuações (scores) e rotulações (labels). Pontuações, quando se atribui uma pontuação a cada instância. Essa pontuação representa o grau de anomalia daquele objeto. No entanto, Rótulos, define uma rotulação categórica para cada instância como, por exemplo, Normal ou Anomalia, Figura 8.

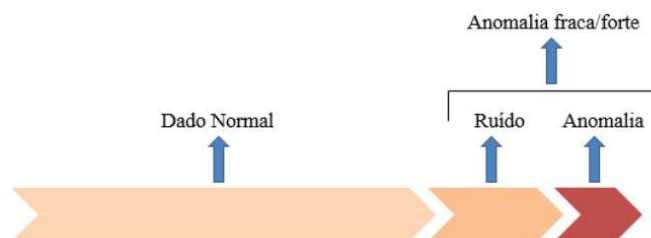


Figura 8: Dados Normais, Ruídos e Anômalos (Fonte: Vasconcelos, 2017).

A seguir serão discutidos os métodos utilizados para detectar anomalias nas séries temporais apresentadas neste trabalho.

3 METODOLOGIAS

*De forma a atender os objetivos proposto neste trabalho, são apresentados a seguir modelos distintos de Detecção de Anomalias, com abordagem Estatística e de Aprendizagem de Máquinas, Figura 9, que foram adotados no desenvolvimento do modelo descrito no item 6.

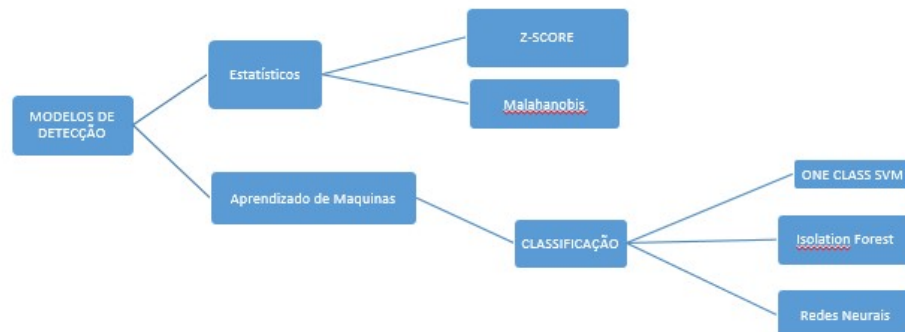


Figura 9: Modelos de Detecção de Anomalias adotados na Investigação.

3.1 MÉTODOS DE DETECÇÃO DE ANOMALIA BASEADOS EM ESTATÍSTICA

Os métodos baseados em estatística assumem que dados normais ocorrem com elevada probabilidade em regiões de um modelo estocástico, enquanto que anomalias ocorrem em regiões com baixa probabilidade (Chandola; Banerjee; Kumar,2009).

Trata-se de uma abordagem mais simples e baseia-se em rotular instâncias que se desviam estatisticamente de propriedades comuns a distribuições: Moda, Média, Mediana e Quartis. Estes métodos podem ser divididos em técnicas paramétricas e não paramétricas:

- a) Paramétricas, que assume uma distribuição prévia estimando parâmetros a partir do conjunto de dados;
- b) Não paramétricas, que não assume nenhuma distribuição previamente.

As abordagens baseadas em estatística possuem como vantagem o fato de fornecerem soluções justificáveis para a detecção de anomalia quando as hipóteses assumidas forem verdadeiras, que combinados com a utilização de uma pontuação (score) de anomalia fornece uma métrica que vai além de somente rotular dados como normais ou anômalos. O score também pode ser usado como intervalo de confiança de uma determinada instância. Além disso, quando se consegue um modelo de estimativa robusto para os dados, as técnicas baseadas em estatística podem ser utilizadas também em dados não supervisionados. No entanto, esta abordagem tem o pressuposto de que os dados são gerados a partir de uma determinada distribuição, o que normalmente nem sempre ocorre.

A seguir são apresentados conceitos e modelos estatísticos adotados no algoritmo deste trabalho.

3.1.1 Z SCORE

O algoritmo Standard Score ou **Z-Score** é uma técnica estatística que possibilita identificar anomalias em dados unidimensionais com uma única análise sobre o fluxo dados. O Z-Score torna diferentes tipos de dados comparáveis e de fácil interpretação (Heiman, 2006). O Z-Score descreve a localização da pontuação (raw score's location) de uma instância de dado (bruto) em termos de quão longe acima ou abaixo da média está medido em unidades do desvio padrão (Heiman, 2006). Um Z-Score igual a zero significa que a instância de dado bruto é igual a média. O Z-Score é calculado como mostrado na equação (1), cujo resultado é adimensional:

$$Z = \frac{X - \mu}{\sigma_x}$$

Eq(1)

Onde:

- Z é o Z-Score de uma instância;
- X, valor da instância;
- μ , média da amostra;
- σ_x , desvio padrão da média.

Desta forma, o Z-Score mede a disparidade em número de unidades de desvio padrão e, por conseguinte, viabilizando a comparações dos dados.

O Z-Score possui dois componentes: sinal, positivo ou negativo, indicando se a pontuação bruta está acima ou abaixo da média, e o valor absoluto Z-Score, indicando a distância em relação à média.

Segundo (Chandola, Banerjee e Kumar, 2019), a regra geral considera todas as instâncias de dado cujo módulo do Z-Score é maior que 3 como anomalias, pois na distribuição Gaussiana a região compreendida entre $[\mu \pm 3\sigma_x]$ contém 99.7% das instâncias, conforme ilustrado na Figura 10:

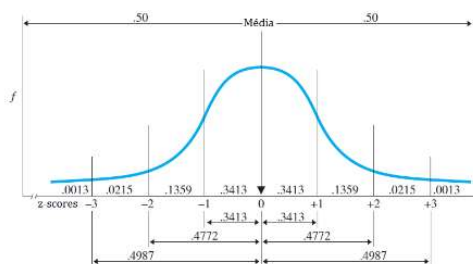


Figura 10: Frequência de Z-Score (Fonte: Heimann, 2006)

Com o valor do Z-Score para cada instância, o algoritmo verifica a distribuição dos Z-Scores (Z-Distribution), isto é, a frequência relativa dos scores brutos de uma população ou amostra. Conforme ilustra a Distribuição Normal de Z na Figura 11, é possível observar que 50% dos scores estão abaixo da média, 50% estão acima da média, e aproximadamente 68% da distribuição está entre $\pm 1 (\sigma_x)$ a partir da média e menos de 1% da distribuição dos scores são maiores que a parcela $\{3\pm(\sigma_x)\}$. Logo, a vantagem do Z-Score é não requerer parâmetros de configuração e assim as anomalias são descobertas conforme a obtenção dos dados. No entanto, o Z-Score pode ter sua capacidade de identificar anomalia com precisão comprometida caso haja poucas instâncias disponíveis (Hodge e Austin, 2004).

Porém, quando os valores da Média, μ , e do Desvio Padrão, σ , são afetadas pelo conjunto de dados, onde 50% dos dados são anômalos, adota-se o Z-SCORE Robusto, onde

valores de Média e Desvio Padrão são substituídos respectivamente, pela Mediana e Mediana do Desvio Absoluto (MAD), conforme equações (2) e (3):

$$MAD = mediana(|X - mediana(x)|) \quad \text{Eq (2)}$$

$$f(X) = \frac{0,6754 * (|X - mediana(X)|)}{MAD} \quad \text{Eq(3)}$$

3.1.2 MAHALANOBIS

Ao se considerar as séries multivariáveis, as anomalias de processos são determinadas por meios de testes de hipóteses que levam em consideração as correlações e colinearidades entre as variáveis medidas. Neste trabalho, utilizar-se-á a técnica Análise de Componentes Principais ou **PCA** (Principal Component Analysis), que determina a anomalia do processo com base na distância dos pontos medidos em relação a uma “média geral” de todas as medições (Junior et al, 2007).

A técnica PCA calcula autovetores e autovalores da matriz de covariância obtida a partir dos dados normalizados. Os autovetores, especialmente aqueles com maiores autovalores associados, fornecem importantes informações sobre o padrão de distribuição dos dados (Hart e Duda, 2000). Desse modo, ordenando os autovalores é possível deduzir a quantidade de informação descrita por cada autovetor. Denominando os autovetores de “componentes”, aqueles associados aos autovalores mais elevados são então chamados “componentes principais”. Os componentes principais que compõem a matriz de componentes (MC).

Assim, o modelo PCA é representado pela matriz “C”, conforme Equação 4. A matriz pode então ser aplicada diretamente aos dados experimentais ($\hat{x}exp_k^T$) obtendo-se os dados modelados ($\hat{x}mod_k$), de acordo com a Equação 5:

$$C = MC.MC^T \quad \text{Eq.(4)}$$

$$\hat{x}mod_k = C.\hat{x}exp_k^T \quad \text{Eq.(5)}$$

Assim, conforme Han, Pei e Kamber (2011), os métodos de Detecção de Outliers univariados podem ser modificados para tratar dados multivariados. Desta forma, pode-se adotar a distância **Mahalanobis** como métrica para análise de séries de dados multivariados do PCA, que é uma medida de distância que considera o formato da distribuição de dados. Para tanto, define-se um limiar para separar dados anômalos. A equação (6) apresenta como se calcula a distância de Mahalanobis entre um ponto X e a média dos dados \bar{X} :

$$mahalanobis(X, \bar{X}) = (X - \bar{X})S^{-1}(X - \bar{X})^T \quad \text{Eq.(6)}$$

Onde S é uma matriz de covariância dos dados. Assim, a distância Mahalanobis de um objeto até a média da distribuição, e está diretamente relacionada com a probabilidade deste objeto. Desta forma, a distância Mahalanobis é igual ao log da densidade da probabilidade do objeto com maior frequência (Tan; Steinbach; Kumar, 2009).

Desta forma, os componentes obtidos no PCA podem ser utilizados com classificadores par identificar as anomalias utilizando a distância Mahalanobis, conforme exemplo da Figura 10:

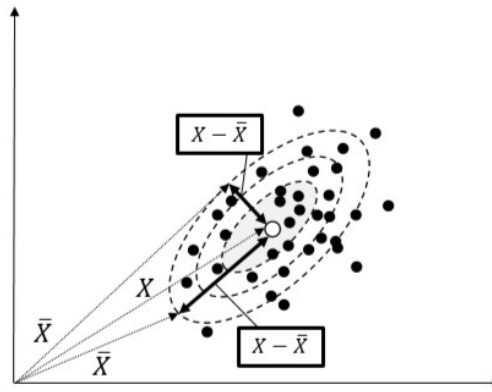


Figura 10: Ilustração de distância de Mahalanobis aplicada a PC.(Fonte: Debes; Koenig e Gross,2005)

Então, a variável calculada da distância de Mahanobolis, (X, \bar{X}) , é uma variável univariada, onde teste de Grubbs pode ser aplicado a esta medida, viabilizando a detecção de anomalia em dados de multivariáveis, executando os seguintes passos:

- i). Calcule o vetor médio do conjunto de dados multivariados;
- ii) Para cada objeto X , calcule Mahanobolis (X, \bar{X}) ;
- iii) Detectar outliers no conjunto de dados univariados transformados, com $Mahanobolis(X, \bar{X}) | X \in D$;
- iv) Se Mahanobolis (X, \bar{X}) indicar como um dado anômalo, então X também é considerado um outlier.

3.2 MÉTODOS DE DETECÇÃO DE ANOMALIA COM APRENDIZAGEM DE MÁQUINA

A seguir são destacados e descritos métodos baseados em técnicas de aprendizagem de máquina para detecção de anomalia baseados em classificação.

3.2.1 MÁQUINA DE VETORES DE SUPORTE DE CLASSE ÚNICA (ONE CLASS SVM)

Uma Máquina de Vetores de Suporte (SVM) introduzido por (Vapnik, 1995), normalmente utilizada para problemas de classificação, foi adaptada para o problema de detecção de anomalias. A extensão One-class SVM proposta por (Schölkopf et al., 2001) é utilizada para detecção semi-supervisionada e não supervisionada, mas em problemas de classificação semi-supervisionada de classe-única (One-Class), as técnicas de detecção de anomalia assumem que são fornecidas todas as instâncias de treinamento com apenas uma única classe, conforme pode ser visto na Figura 11. Assim as fronteiras discriminativas são aprendidas em torno dessa classe com aplicação de Kernels robusto sem regiões mais complexas.

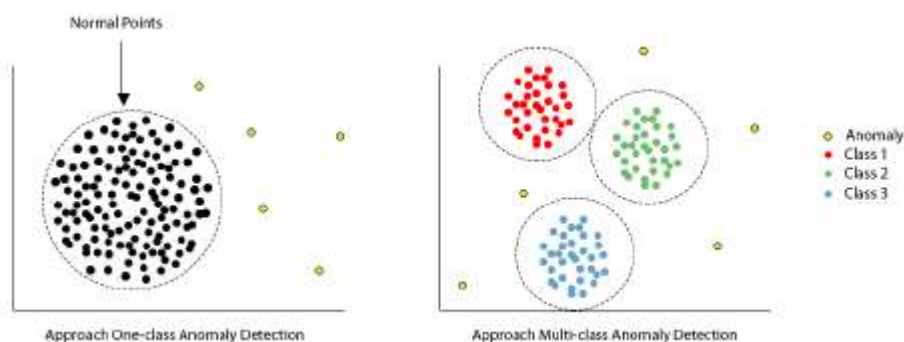


Figura 11: Classificação conforme One Class Support Vector Machine (Fonte: Mata, 2017).

Assim, o algoritmo é capaz de aprender uma fronteira suave para agrupar instâncias normais usando o conjunto de treinamento e, em seguida, no conjunto de testes as instâncias são identificadas quando estão fora da região aprendida. Dependendo do caso, o output pode ser numérico ou textual. No cenário não supervisionado One-class SVM é treinado utilizando um conjunto de dados e, em seguida, atribui-se um score para cada instância do conjunto de dados por meio do uso de uma distância normalizada determinada pela fronteira de decisão (Amer, Goldstein; Abdennadher, 2013)

3.2.2 ISOLATION FOREST

Isolation Forest (Liu; Ting; Zhou, 2008), é um modelo não supervisionado baseado em 'floresta de decisão' que tem por objetivo detectar anomalias presentes num conjunto de dados a partir dos atributos fornecidos. Uma vez que, em geral, anomalias ocorrem em uma proporção bem menor do que os dados normais e, além disso, apresentam valores de atributos muito distintos, criar ramificações baseadas em valores aleatórios dos atributos (escolhidos entre os valores mínimo e máximo presentes na base de treino do atributo em questão) pode separar melhor as anomalias, já que tenderão a se localizar nos primeiros nós das árvores. A Figura 12 extraída de Liu, Ting e Zhou (2008), ilustra a separação de anomalias por partições aleatórias,

onde perceber que as anomalias precisam de uma quantidade menor de ramificações para serem isoladas.

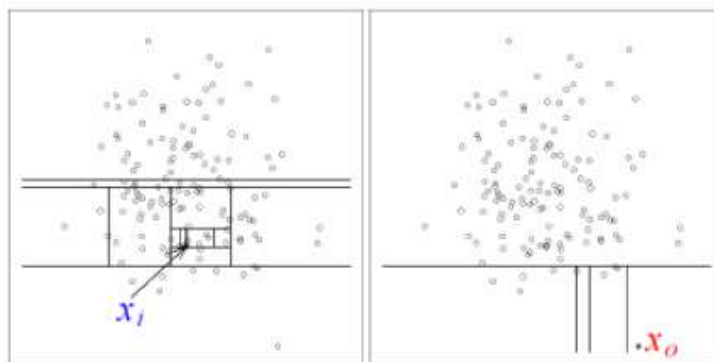


Figura 12: Ilustração do processo de isolamento de anomalias x_0 entre amostras normais x_i . (Fonte: Liu; Ting; Zhou, 2008).

Após a criação das árvores, novas anomalias são detectadas utilizando uma proporção entre a esperança do número de vértices percorridos por uma entrada, X_e , o número médio de vértices a se percorrer em um conjunto de dados de tamanho N . A partir dessa métrica, definida em Liu, Ting e Zhou (2008), pode-se concluir que se o número de vértices percorridos for próximo a 0, então a entrada fornecida é uma anomalia, caso contrário; a entrada pode ser considerada como normal

3.2.3 REDE NEURAL ANTOENCODER

As Redes Neurais são compostas por neurônios artificiais, que são unidades básicas que realizam cálculos simples, e interagem com outros neurônios sob diferentes formas e organizações. Sua arquitetura pode ser do tipo:

- a) FEEDFORWARD, com conexões acíclicas entre neurônios da mesma camada, mais aplicada a problemas Supervisionados, de classificação e regressão. ou;
- b) RECORRENTES, com conexões cíclicas entre neurônios, aplicadas a previsões.

As redes FeedForward, Figura 13, adotada neste trabalho, são caracterizadas pelo arranjo dos neurônios em camadas, cujos neurônios ou Perceptron da mesma camada não se conectam entre si, mas conectam com neurônios da camada subsequente (estado Forward) ou anterior, durante o processo de aprendizagem pela atualização de “pesos” no cálculo, (estado Backward).

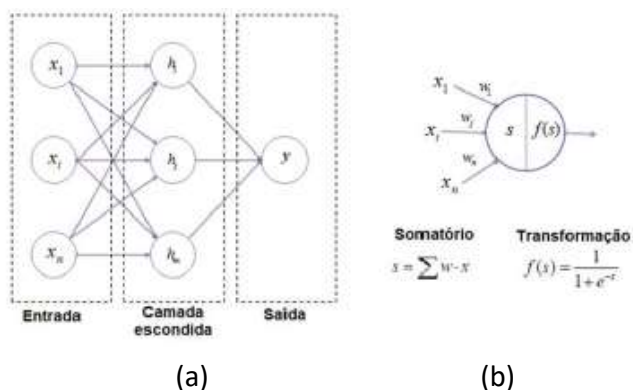


Figura 13: a) Rede Neural do tipo MLP, b) Perceptron (Sayad, 2019).

A rede Multi Layer Perceptron, MLP, mais utilizada na literatura, proposta por Rumelhart e Hilton (1986) contém três partes principais: a camada de entrada, a camada de saída e um conjunto de camadas ocultas, que pode conter várias camadas vinculadas entre elas. O uso de várias camadas ocultas permite que a rede aprenda padrões e dependências complexas entre os dados de entrada e saída.

Todas as camadas são ligadas por parâmetros chamados ‘pesos’, representados pelas linhas na Figura 13, e para cada neurônio, representado pelos círculos na Figura 13, existe uma função de ativação f que decide se o estado do neurônio mais próximo estará ligado (estado alto) ou desligado (estado baixo). No caso de redes MLP, essa função de ativação é geralmente um sigmoide, apropriada para a classificação e a escolha, pois limitam o resultado entre os valores 0 (desativado) e 1 (ativado).

Desta forma, o objetivo dessa rede neural é prever os mesmos valores de saída esperados usando os valores de entrada. Isso é feito alterando os parâmetros de peso que vinculam suas camadas. Para isso, uma etapa de aprendizado é necessária, em que esses parâmetros serão otimizados de uma maneira controlada para obter o melhor resultado no final.

O resultado de saída durante o treinamento da rede neural é comparado com a saída real desejada, através de uma função erro, e assim por meio do método iterativo Backpropagation, os parâmetros da rede são reajustados com base no valor da saída e propagados para o início da rede, na primeira camada, estado de aprendizado Backward (Rumelhart e Hilton, 1986).

Conforme Lecun, Bengio e Hinton (1986), a otimização dos parâmetros no Backpropagation utiliza um algoritmo de gradiente descendente, que calcula derivadas parciais da função de cálculo dos erros residuais entre saída esperada e obtida, alterando o peso nos neurônios da camada inicial e nas demais das camadas subsequentes a cada ciclo de aprendizagem, resultando em um efeito propagado nas camadas anteriores de forma que os pesos dos neurônios sejam ajustados para aproximação do resultado esperado. A cada iteração o gradiente descendente é calculado de forma a reduzir o erro, para tornar a derivada da função erro nula, ou seja, na região do erro mínimo local.

Entre as redes Neural Feedforward, destaca-se a Autoencoder (AE), um modelo de Deep learning, conforme Makhzani, Frey e Goodfellow, (2014), cujos dados de entrada e saída, são submetidos respectivamente a fase de compressão (Encoder) e descompressão (Decoder). Esta arquitetura utiliza o Autoencoder na tarefa de redução de dimensionalidade (DR), aprendizado de parâmetros durante a compressão e de modelos generativos. Assim, o Autoencoder aprende características inerentes dos dados, criando um espaço dedicado dimensional reduzido, que descarta dados redundantes, não representativos e ruídos, Figura 14.

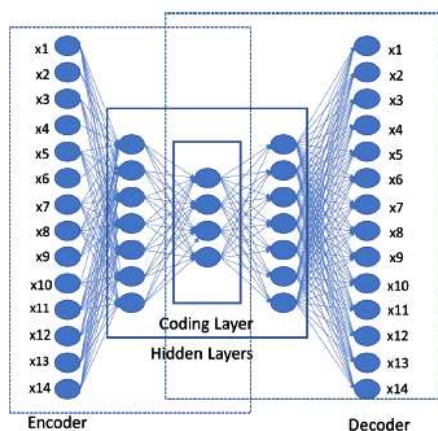


Figura 14: Rede Neural Autoencoder.

3.4 MÉTRICAS DE DESEMPENHO

Embora haja distintas métricas de desempenho existentes na literatura, neste trabalho, são utilizadas as métricas aplicadas em Vargas; Munaro; Ciarelli (2016) conforme definições descritas a seguir utilizadas na Matriz de Confusão obtida com os resultados, Figura 15:

- Verdadeiro-positivo (*VP*): número de amostras classificadas normais corretamente;
- Falso-negativo (*FN*): número de amostras classificadas normais equivocadamente;
- Verdadeiro-negativo (*VN*): número de amostras anômalas classificadas corretamente;
- Falso-positivo (*FP*): número de amostras anômalas classificadas equivocadamente;

		VALOR PREDITO	
		SIM	NÃO
REAL	SIM	VP	FN
	NÃO	FP	VN

Figura 15: Matriz de Confusão

A partir dos valores da Matriz de Confusão, *VP*, *FN*, *FP* e *VN*, são calculadas as seguintes métricas:

- Acurácia (*A*): Razão entre soma de Valores Verdadeiros e quantidade de dados;
- Precisão (*P*): razão entre Verdadeiros Positivos (*VP*) e os Resultados Positivos (*VP+FP*), vide Equação (7);
- Sensibilidade (*S*): razão entre Verdadeiros Positivos (*VP*) e a soma dos Verdadeiros Positivos com Falso Negativos (*VP+FN*), vide Equação (8);
- Medida F (*F1*): média harmônica entre precisão e sensibilidade, vide Equação (9):

$$P = \frac{VP}{VP+FP}$$

Eq.(7)

$$S = \frac{VP}{VP + FN}$$

Eq.(8)

$$F1 = \frac{P}{S}$$

Eq.(9)

Neste trabalho, as avaliações cujo foco estava na detecção de anomalias foram classificadas com a métrica precisão.

4 ARQUITETURA DO SISTEMA PROPOSTO

No trabalho em questão, seguindo a linha de desenvolvimento de autores como Dominguesa et al (2018), adotou-se a estratégia de desenvolvimento de um algoritmo que teste distintos métodos de Detecção de Anomalias com distintas abordagens (estatística e de aprendizagem), elencadas e agrupadas na Tabela 1, discutidas anteriormente no Capítulo 3 deste trabalho: considerando as naturezas diferentes de processo e características de construção dos sistemas de compressão de diversas plantas offshore para assim viabilizar escalabilidade do seu uso conforme resultados históricos de cada planta, ou seja fornecem domínios distintos (envelope operacional ou domínio).

Tabela 1: Métodos e bibliotecas Python utilizadas.

MÉTODO	BASEADO EM	APRENDIZADO DE MÁQUINA	COMANDOS (BIBLIOTECAPYTHON)
Z-SCORE	ESTATISTICA	NÃO	SCIPY. STATS
MAHALANOBIS	ESTATISTICA MUTIVARIÁVEL	NÃO	PCA (SKLEARN. DECOMPOSITION) StandardScaler (SKLEARN.PREPROCESSING) Empirical Covariance, MinCovDet(SKLEARN.COVARANCE)
OC-SVM	CLASSIFICAÇÃO OU DENSIDADE	SUPERVISIONADO	OneClass SVM (SKLEARN. SVM)
ISOLATION FOREST	CLASSIFICAÇÃO OU DENSIDADE	NÃO SUPERVISIONADO	Isolation Forest (SKLEARN. ENSEMBLE)
REDES NEURAL AUTOENCODER	CLASSIFICAÇÃO	NÃO SUPERVISIONADO	Sequential, Model, Input, Conv2D, Flatten, Dense, Regularizes, Plot Model(Tensorflow. Keras/.Layers/.Utils)

Assim, o modelo foi desenvolvido em Python, conforme algoritmo descrito na Figura 16 a seguir:

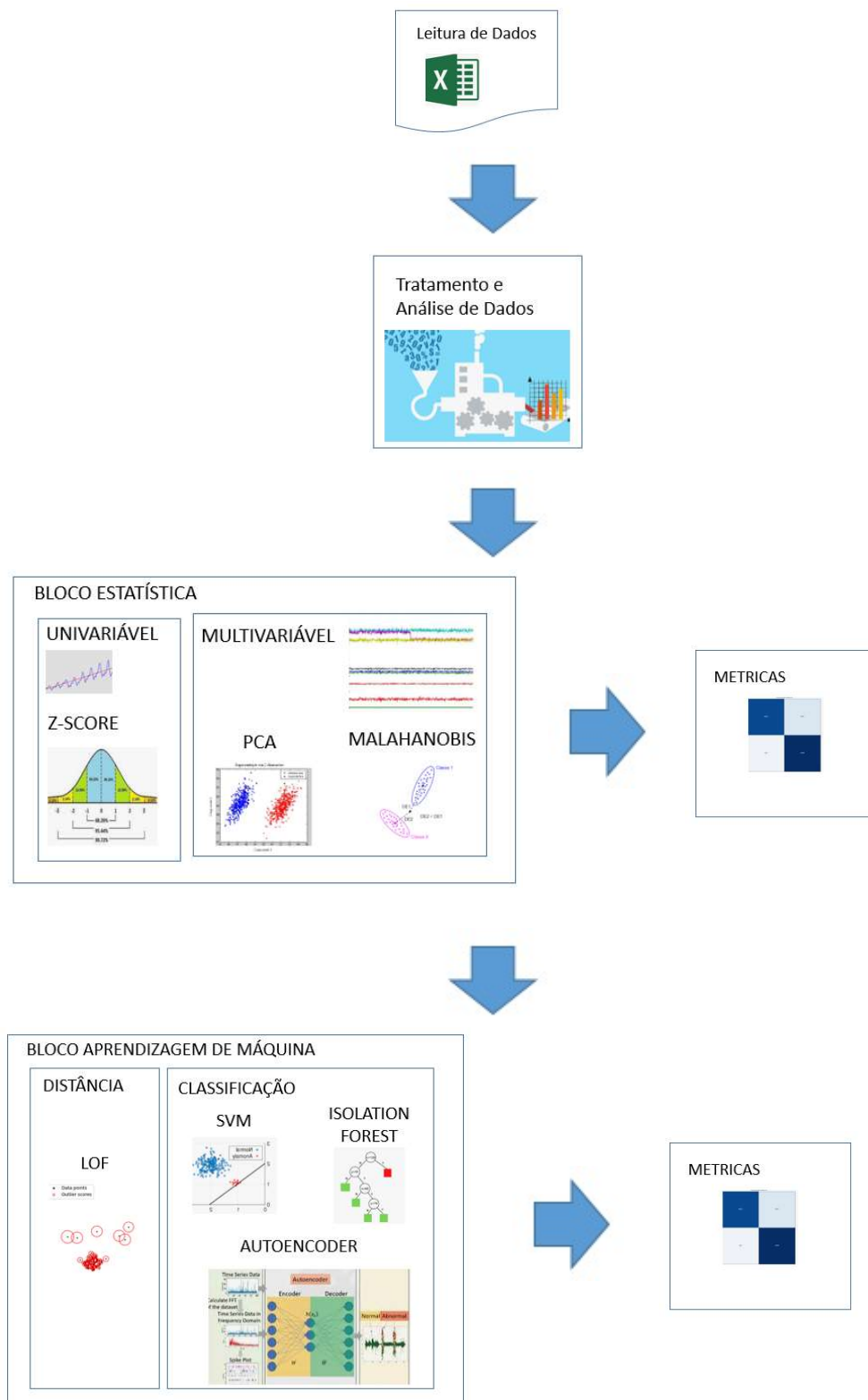


Figura 16: Algoritmo de Detecção de Anomalias

Para definir base de dados na implantação dos métodos de Detecção de Anomalia, em primeiro lugar foram identificadas em um banco de dados de parada do sistema de compressão adotado na investigação, os sinais de sensores que mais contribuíram para iniciar o processo de intertravamento ou desligamento que leva a parada dos equipamentos, vide gráfico da Figura 17:

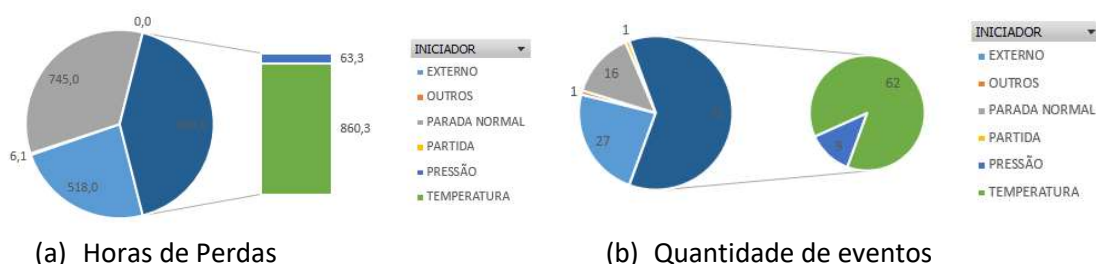


Figura 17: Pareto de Eventos Iniciadores em Perdas de Produção de Sistema de Compressão, (a) Horas, (b) Quantidade de Eventos

Desta forma, os sensores de temperatura e pressão destacaram-se como iniciadores ou motivadores de paradas não forçadas. Desta forma foram identificadas as datas e horários destas paradas em 2017 e 2018. Considerando estas datas foram extraídas séries temporais de dados de temperatura, pressão e inclusive vazão, Figura 18, do circuito de compressão investigado, em períodos de até 1 mês que antecederam os eventos de parada, a partir do servidor do PI Processbook (OSYSOFT), vide arquitetura de extração e tratamento de dados (cerca de 43000 linhas), na Figura 19:

	Data	Temperatura_Entrada	Temperatura_Saida	Pressao_Entrada	Pressao_Saida	Vazao	Anomalia
0	15/05/2017 00:00	67,5702057	169,8499146	61,96791077	21,88606644	60838,55078	0
1	15/05/2017 00:01	67,5946274	169,4202728	61,41078568	21,78266335	61253,11328	0
2	15/05/2017 00:02	67,5900803	168,9943085	61,02445602	21,7167778	61544,19141	0
3	15/05/2017 00:03	67,5611191	168,6232605	60,53781509	21,62524033	62890,96875	0
4	15/05/2017 00:04	67,4969482	168,0438232	60,04884338	21,59707069	63825,34766	0

Figura 18: Trecho dos dados do circuito de compressão investigado.

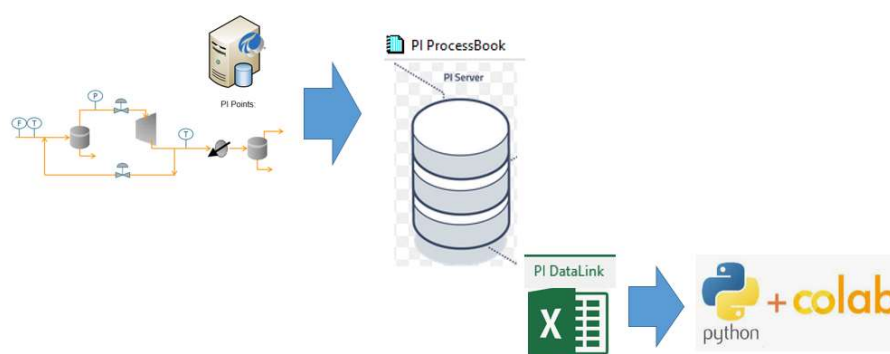


Figura 19: Arquitetura de Extração e Tratamento de Dados

Deve-se destacar que para efeito de preservar a confidencialidade dos dados de processo da empresa, informações acerca da identificação da planta, seus sistemas e sensores foram então descaracterizados.

Os dados de processos gravados no servidor do PI Processbook, então foram submetidos a cálculo de média minuto a minuto, e extraídos através do PI DataLink que é um suplemento do MS Excel que permite recuperar informações do PI Server diretamente. para uma planilha, convertida em arquivo extensão para ser consumido pelo Notebook Python desenvolvido no Google Colab, ANEXO TRABALHO FINAL-PUC. Jpynb.

5 RESULTADOS

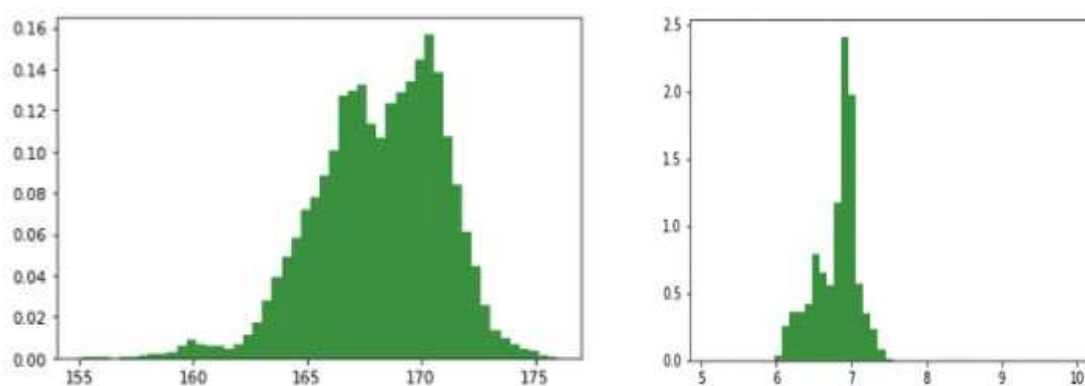
No Notebook desenvolvido no Python, conforme primeiros blocos da Figura 16, foram implementadas as bibliotecas necessárias a coleta, tratamento, como conversões de tipo de dados, vide trecho da tabela de dados, eliminação de *missing values*, e definição das restrições de domínio, e análise exploratória dos dados

Então a partir da extração de dados, foi realizada a análise exploratória dos dados, com cálculo de variáveis estatísticas (média, desvio, mínimo, ...) para a base de dados de cada variável coletada (temperaturas, pressões e vazões), conforme exemplo da Figura 20:

```
Temperatura Saída
count    42726.000000
mean      168.168638
std        2.805396
min       150.013534
25%       166.381256
50%       168.425156
75%       170.288055
max       175.981247
Name: Temperatura_Saída_Convertido, dtype: float64
```

Figura 20: Dados estatísticos da variável temperatura de saída.

Em seguida foram elaboração dos histogramas dos dados para avaliar sua distribuição conforme exemplo da Figura 21:



a) Temperatura

b) Pressão

Figura 21: Histograma dos Dados, a) Temperatura, b) Pressão;

Além disso, verificou-se a Normalidade dos dados, através do gráfico QQ, conforme a Figura 22:

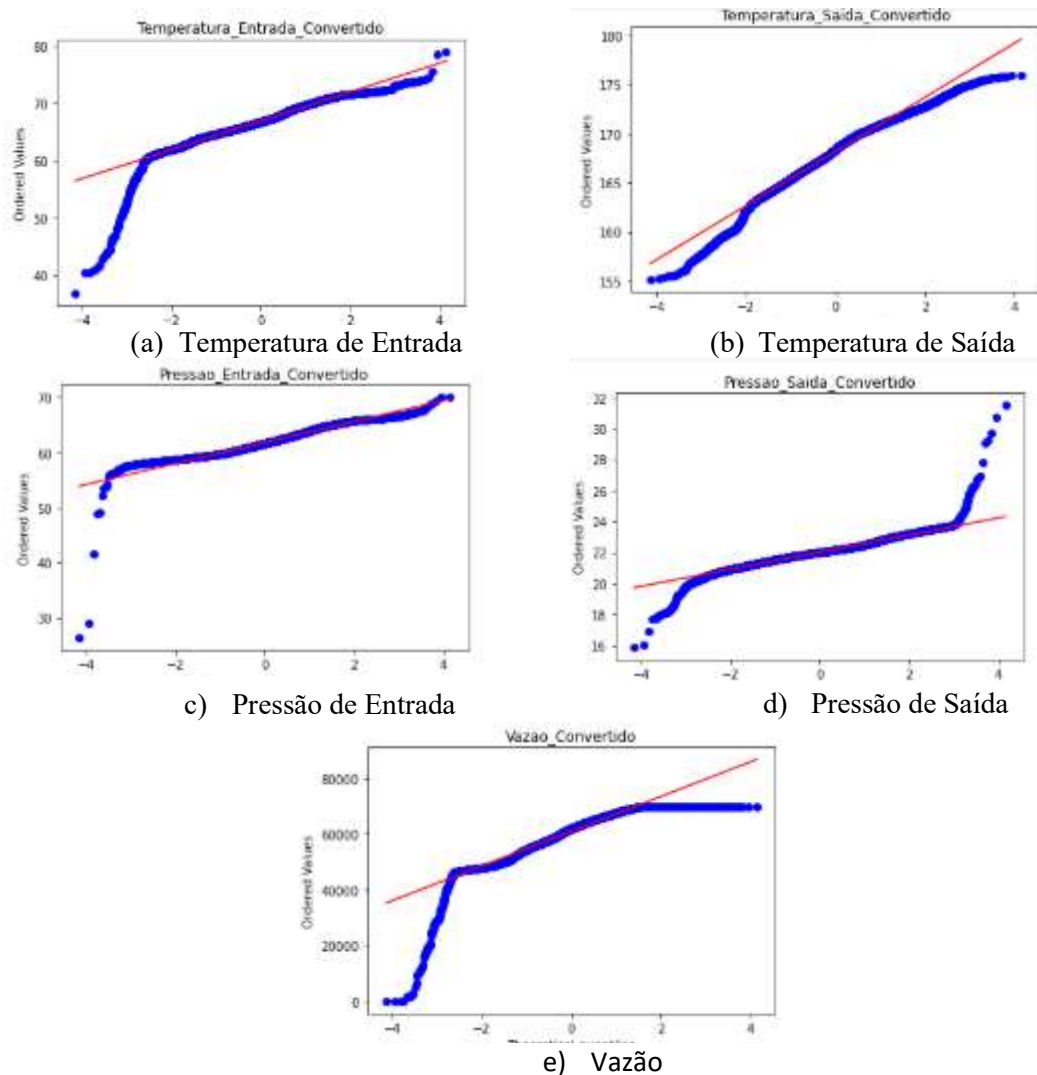


Figura 22: Análise Gráfica QQ dos dados coletados para cada variável.

Durante a análise dos gráficos QQ, identificou-se que a variável vazão tem um limite de range de medição restrito a 70000 m³/h, o que promove a existência de valores faltantes para quando a vazão real estiver acima deste valor. Desta forma, a variável vazão não foi considerada como entrada na análise multivariável de detecção de anomalia deste trabalho, devido ao limite operacional do instrumento.

Com base na análise de normalidade os envelopes operacionais ou domínio foram ajustados, considerando somente valores acima dos limites inferiores em que os sistemas estavam operacionais, uma vez que objetivo na detecção de anomalias durante operação, gerando um novo conjunto de dados, agora com tratamento de domínio, adotados para a continuidade da aplicação do Notebook, cuja distribuição é apresentada na Figura 23, assim como a análise gráfica de normalidade QQ, apresentada na Figura 24:

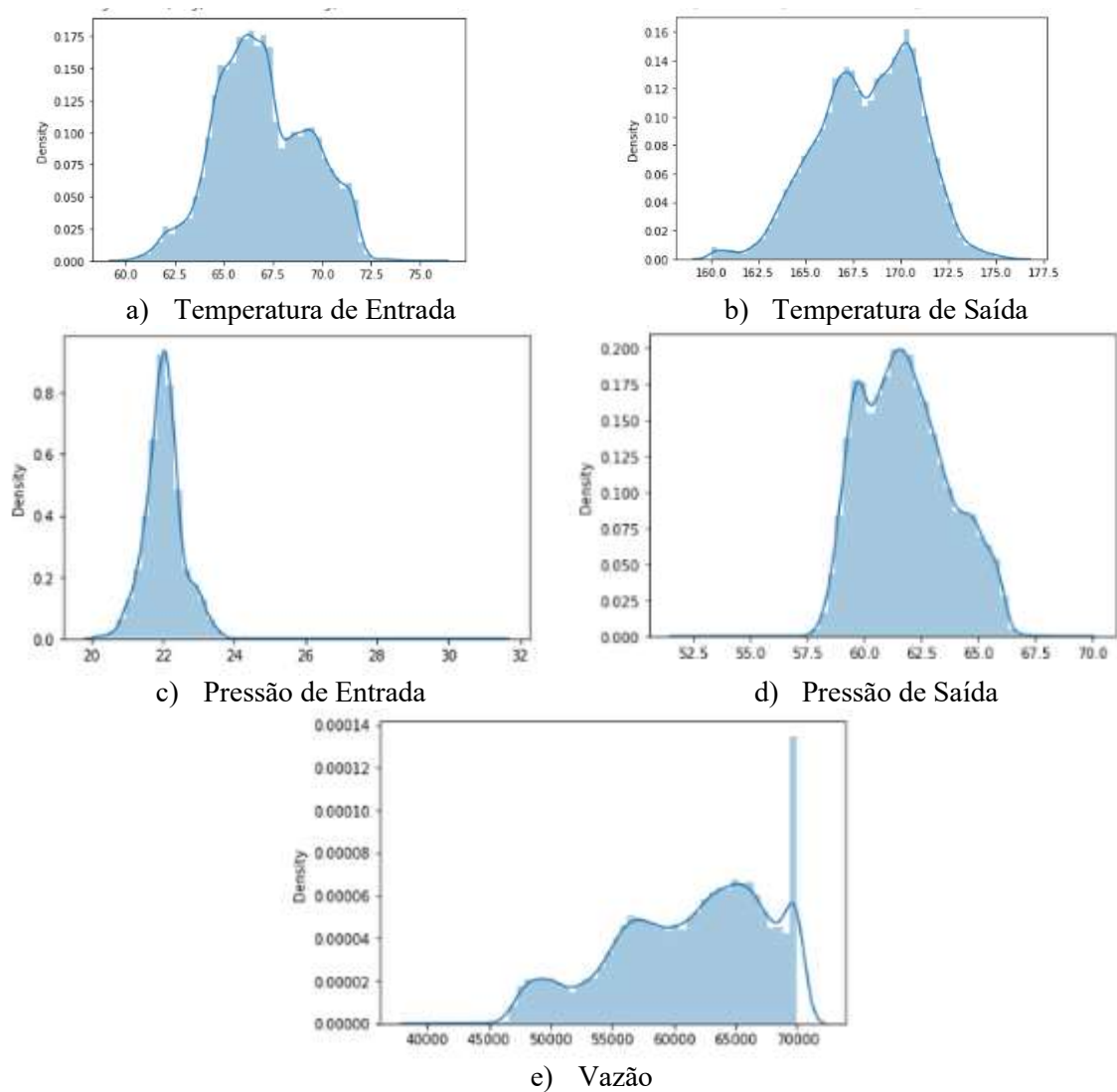


Figura 23: Distribuição de Dados com tratamento de domínio.

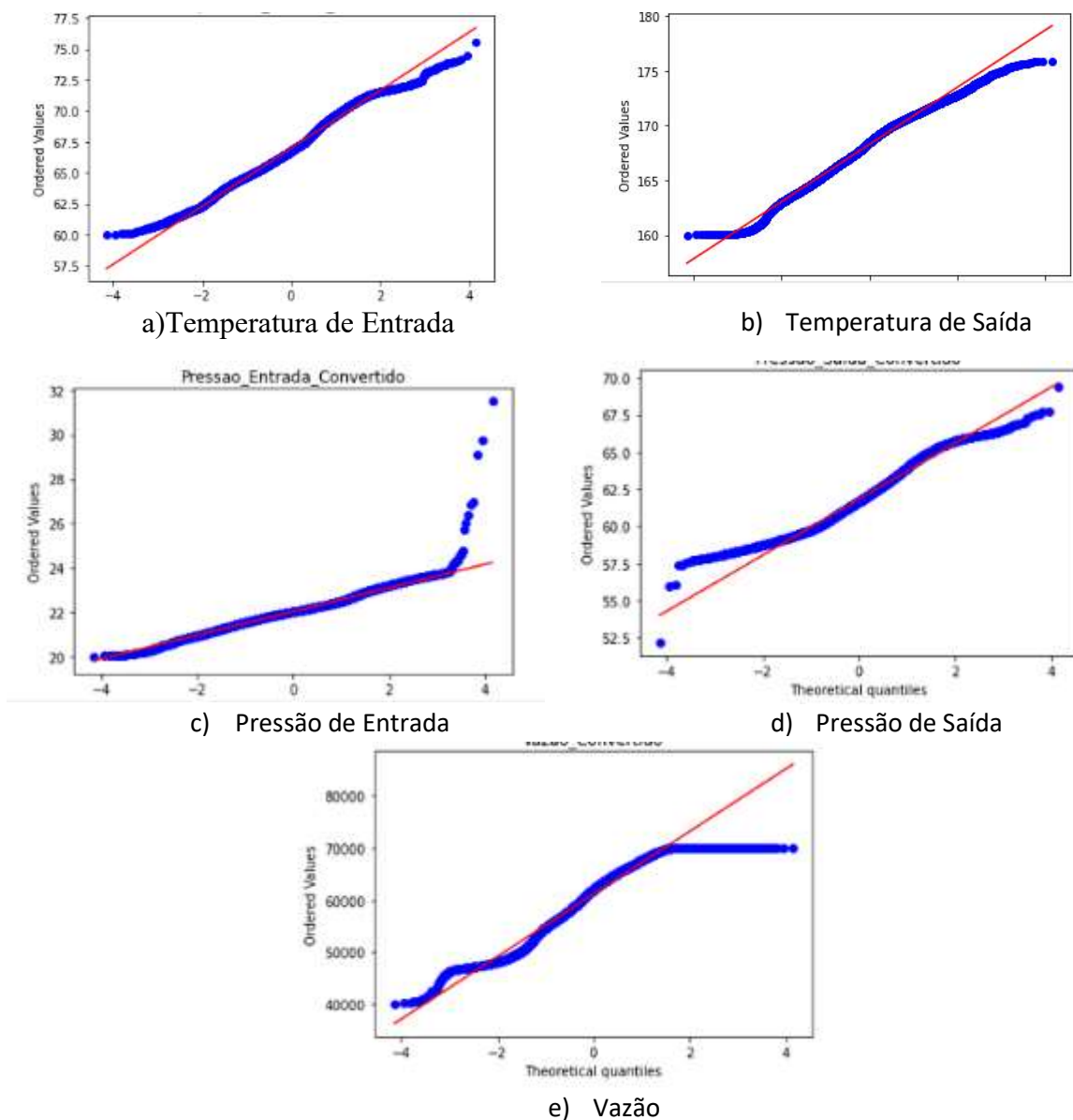


Figura 24: Análise Gráfica QQ com tratamento de domínio.

Para finalizar a análise exploratória dos dados, foi realizada a análise da correlação entre as variáveis, seja quantitativa conforme a Figura 25, com uso do comando Heatmap, e através da análise gráfica, vide Figura 26, através do comando Scatter Plot:

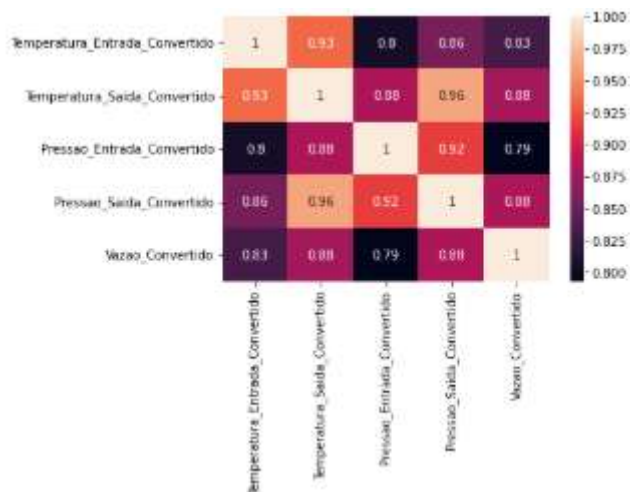


Figura 25: Matriz de Correlação das variáveis de temperatura e pressão com ferramenta Heatmap.

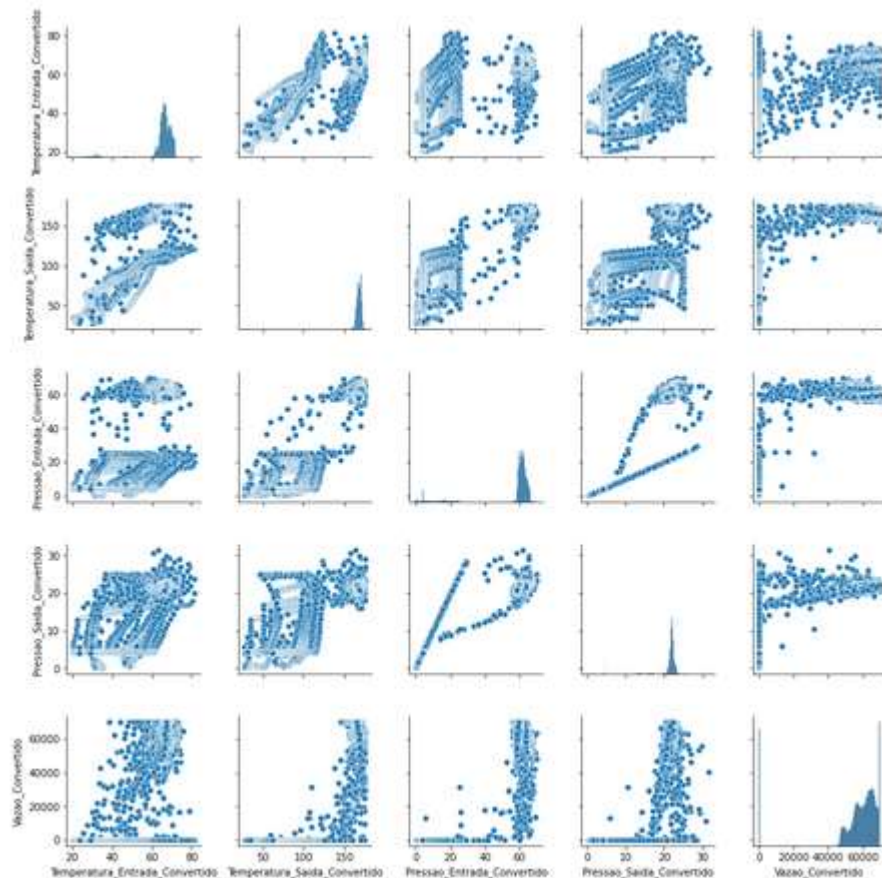


Figura 26: Análise da correlação dos dados de temperatura e pressão através do comando Scatterplot.

Observou-se uma correlação maior entre pressão de saída e temperatura de saída, de 0,96, e menor, de 0,79 da vazão com a pressão de entrada, ambas justificadas justificada pela relação física que rege o processo de compressão isentrópica, modelada pela Lei dos Gases Perfeitos, vide Equação 10:

$$(\text{Pressão} \cdot \text{Volume}^k = nRT \text{ Temperatura}). \text{Eq}(10)$$

A análise de correlações entre as distintas variáveis, demonstraram viável a redução de dimensionalidade na análise de multivariáveis. Os dados foram separados para realizar o aprendizado de máquinas de Treino com dados de 2017 e de Teste, com dados de 2018.

De posse das variáveis estatísticas, tais com média, desvio padrão dos dados, foi possível conduzir o primeiro modelo de Detecção de Anomalias Z-SCORE, conforme Figura 27:

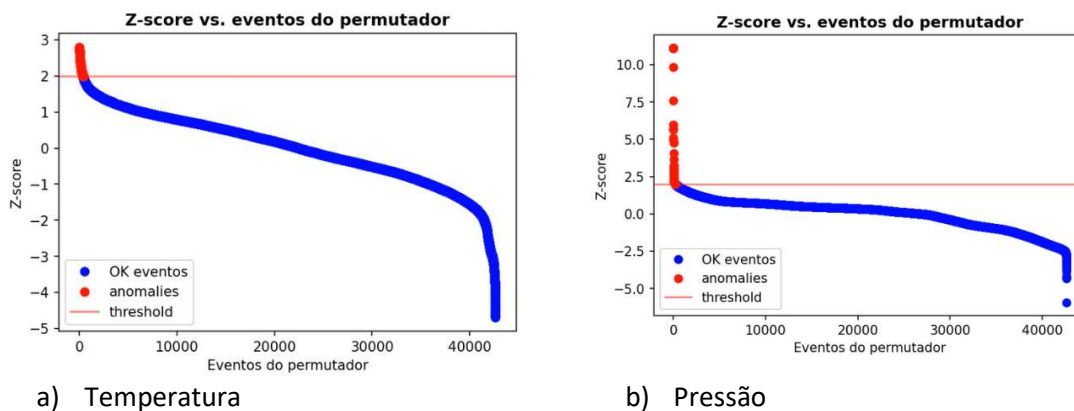


Figura 27: Z-Score, a) Temperatura, b) Pressão;

Em seguida, ainda, foram adicionadas a abordagem estatística para séries multivariáveis com PCA, decompostas, Figura 28, e aplicada a distância de Mahalanobis, desta forma concluindo a sequência de modelos estatísticos de detecção de anomalias, onde valores acima do Treshold (destacado em vermelho) indicam a presença de outliers, Figura 29:

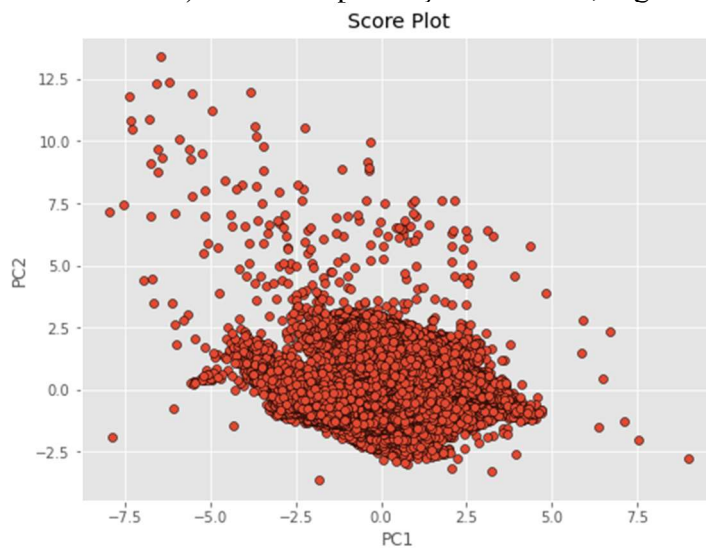


Figura 28: Resultados do PCA

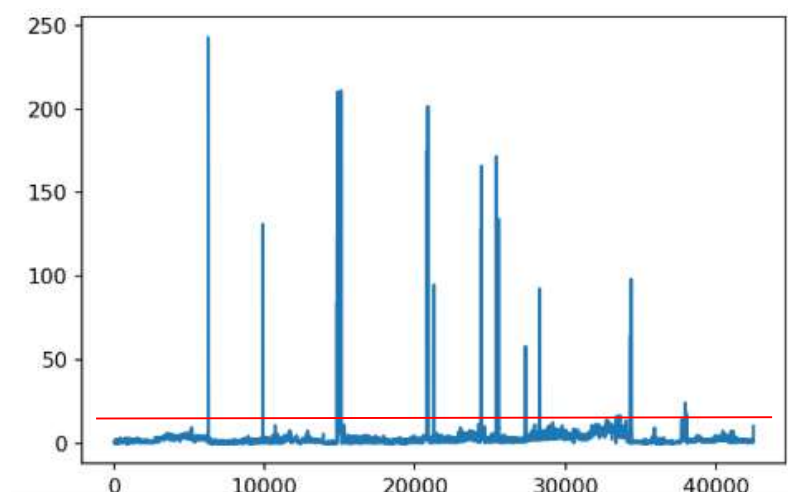


Figura 29: Distância de Mahalanobis calculado para os dados da série multivariada

Dando sequência, seguiu-se a implementação de modelos de aprendizagem, relacionados a classificação: SVM, Isolation Forest e Rede Neural Autoencoder.

Assim, todos os modelos implementados foram submetidos a treino, e com exceção da Rede Neural Autoencoder foram otimizados, conforme parâmetros da Tabela 2: Em todos os blocos foram calculadas métrica de Precisão baseadas na Matriz de Confusão utilizada como alvo no processo de otimização durante o treino.

Tabela 2: Parâmetros de Otimização dos Modelos

MODELO	PARAMETROS
Z-SCORE	Desvio Padrão=2,1; 2,2; 2,3; 2,4; 2,5; 2,6 ; 2,7; 2,8; 2,9; 3,0 Janelas de Média Móvel =1 a 10
MAHALANOBIS	K=1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 30, 40, 50
OC-SVM	Estimadores=100,200,300,400, 500
ISOLATION FOREST	Tipo de Kernel: Linear, Poly, RBF, Sigmoid Nível: 1, 3, 5
REDE NEURAL AUTOENCODER	Otimizador=adadelta e loss=mse

Em seguida os modelos já otimizados foram avaliados com a base de testes, conforme resultados de Precisão apresentados na Tabela 3:

Tabela 3: Resultados de Precisão de Treino e Teste na aplicação dos modelos de Detecção Anomalias.

MODELO	PRECISÃO		PARAMETROS OTIMIZADOS
	TREINO	TESTE	
Z-SCORE	0,5	0,0	Z=2,8 Desvios, Média móvel= 9 pontos
MAHALANOBIS	0,002118	0,0	K=4
OC-SVM	0,00009549	0,00010744	Kernel=RBF,degree=1
ISOLATION FOREST	0,0023	0	300
REDE NEURAL AUTOENCODER	0	0	-

Porém deve-se destacar que os modelos identificaram dados anômalos ou Falso Positivos, conforme Tabela 4, com exceção do modelo Autoencoder, que são úteis como alertas no monitoramento do sistema de compressão, por identificarem cenários de operação fora do regime normal, cuja evolução leva à condições extremas e ao consequente intertravamento do sistema.

Tabela 4: Resultados de Falso Positivo na aplicação dos modelos de Detecção Anomalias.

MODELO	FALSO POSITIVO	
	TREINO	TESTE
Z-SCORE	1	0
MAHALANOBIS	471	571
OC-SVM	20941	9306
ISOLATION FOREST	1200	304
REDE NEURAL AUTOENCODER	0	0

6 CONCLUSÕES E TRABALHOS FUTUROS

A revisão bibliográfica realizada neste trabalho permitiu identificar modelos de base estatística e de aprendizagem de máquinas com aplicação à detecção de anomalias em séries temporais multivariados, incluindo a comparação entre estes métodos.

A análise exploratório dos dados e seu tratamento evidenciou a correlação entre as distintas séries, de forma que mesmos tivessem sua dimensão reduzida, além de caracterizar a normalidade dos dados.

Foi desenvolvido um algoritmo em Python que submeteu o conjunto de dados de sensores de processo de um sistema de compressão (temperatura e pressão) aos seguintes modelos de Detecção de Anomalias: um primeiro grupo estatístico, Z-Score e Distância de Mahalanobis, e um segundo grupo de aprendizagem de máquinas, baseados em classificação, OC-SVM, Isolation Forest e Rede Neural Autoencoder.

Os resultados da aplicação e otimização destes modelos no conjunto de dados investigado identificou que:

- O processo de otimização dos hiperparâmetros permitiu alavancar os resultados de precisão nos modelos, com exceção no modelo de Autoencoder que não localizou Falso Positivos;
- Na base de treino com dados de 2017, os modelos apresentaram resultados de precisão superiores, com destaque para modelo Z-SCORE (Estatístico), seguido do modelo Isolation Forest (Aprendizado de Máquina), demonstrando um descolamento dos dados de teste com dados de 2018, que pode estar associado a um comportamento operacional do sistema de compressão no ano de 2018;
- Na base de teste, somente o modelo OC-SVM apresentou resultados de precisão de mesma ordem de grandeza

Assim, para a base de dados os modelos mais recomendados são, entre estatísticos, Mahalanobis, e de aprendizado de máquinas, Isolation Forest.

Deve-se destacar que o critério de classificação para os dados como anômalos foi intertravamento ou desligamento automático da máquina, executado quando os dados ultrapassam limites superiores ao especificado para o equipamento operar de forma segura. No entanto, todos os modelos identificaram anômalos (Falso Negativos) que podem ser utilizados como alertas no monitoramento do sistema de compressão, de forma que ações mitigadoras sejam desencadeadas para migrar de uma zona operacional de risco para uma condição de regime operacional seguro de forma e evitar o intertravamento do equipamento.

7. REFERÊNCIAS BIBLIOGRÁFICAS

- Amer, M.; Goldstein, M.; Abdennadher, S. Enhancing One-Class Support Vector Machines for Unsupervised Anomaly Detection. Proceedings of the ACM SIGKDD Workshop on Outlier Detection and Description (ODD'13), ACM Press, p. 8–15, 2013.
- Bidgoli, A. A., Simulation and Optimization of Primary Oil and Gas Processing Plant of FPSO Operating in Pre-Salt Oil Field, Polytechnic School of the University of São Paulo, 2018.
- Box, G. E. P. and Jenkins, G. M. (2008). Time series analysis: forecasting and control, 4nd. ed., 2008.
- Brand, T., Demands on Sensors for Future Servicing: Smart Sensors for Condition Monitoring, www.analog.com, 2017
- Breunig, M. M. et al. Lof: Identifying density-based local outliers., 2000.
- Carapeto, L.A.V., Controle da Compressão de Gás em Plataforma Offshore: Camada Regulatória, Politécnica da Universidade Federal do Rio de Janeiro, Projeto de Graduação, 2016
- Chandola, V., Banerjee, A., Kumar, V. “Anomaly detection,” ACM Computing Surveys, vol. 41, no. 3, pp. 1–58, Jul. 2009
- Cheboldi, D., Anomaly Detection of Time Series, Thesis, University of Minnesota, 2010.
- Cheng, H., Tan, P. N., Potter, C., and Klooster, S. (2008). A robust graph-based algorithm for detection and characterization of anomalies in noisy multivariate time series. In 2008 IEEE International Conference on Data Mining Workshops, pages 349–358.
- Debes, K.; Koenig, A.; Gross, H.-M. Transfer functions in artificial neural networks a simulation-based tutorial. *Brains, Minds and Media*, v. 2005, n. 1, 2005.
- Dominguesa, R., Filippone, M., Michiardi, P. Zouaoui, J., A comparative evaluation of outlier detection algorithms: experiments and analyses, San Francisco: Holden-Day, Jihane Elsevier, 2018
- Freitas, I. W. S.. Um estudo comparativo de técnicas de detecção de outliers no contexto de classificação de dados/ Igor Wescley Silva de Freitas. - 2019., Mahalanobis
- Han, J.; Pei, J.; Kamber, M. Data mining concepts and techniques. [S.l.]: Elsevier, 2011.
- Hart, P., Duda, R.O. Pattern Classification. 2ª ed. John Wiley Pro., 2000.
- Heiman, G. W, Basic Statistics for the Behavioral Sciences, 6th ed. Cengage Learning, 2006
- Hodge V. J. and Austin, J., “A Survey of Outlier Detection Methodologies,” *Artificial Intelligence Review*, vol. 22, no. 2, pp. 85–126, Oct. 2004

Janke, P., Machine Learning Approaches for Failure Type Detection and Predictive Maintenance Master Thesis, Technische Universität Darmstadt, 2015.

Junior, C. A. V. Junior, Araujo, O. Q. F., Medeiros, J. L., Classificadores Hierárquicos e PCA na Detecção e Identificação de Falhas em Redes de Escoamento, 4o PDPETRO, Campinas, SP 21-24 de Outubro de 2007

Lecun, Y. Bengio, Y., Hinton, g., Deep Learnig, Nature, v.521, n.7553, p. 436-444, 2015. Disponível em <https://www.nature.com/articles/nature14539>

Liu, F. T.; Ting, K. M.; Zhou, Z. Isolation forest. In: 2008 Eighth IEEE International Conference on Data Mining. [S.l.: s.n.], 2008

Kamal, P., Sugandhi, R., Anomaly Detection for Predictive Maintenance in Industry 4.0 – A Survey, E3S Web of Conferences, 2020.

Makhzani, A., Frey, B. Goodfellow, I. Adversarial Autoencoder, 2014.

Mata, F. D. G., Investigando Métodos Inteligentes para Detecção de Anomalias em Comportamento de Insetos Sociais, Dissertação, Universidade Federal do Pará, 2017.

Rob Hyndman et al. Forecasting Functions for Time Series and Linear Models, 2019

Rumelhardt, D.E., Hilton, G. E., Wilians, R. J., Chapter 8: Learning Internal Representations by Error Propagation. Parallel Distributed Processing: Explorations in the Microstructure of Cognition: Foundations Processing, Vol.1. Foundations. MIT Press. Cambridge, MA., 1986.

Sayad, S., Multilayer Perceptron, disponível em:
http://www.saedsayad.com/artificial_neural_network_bkp.htm

Schölkopf, B. Platt, C. C., Shame-Taylor, J., Smola, A. J., Williamson, R. C., Estimating the support of a high dimensional distribution. Neural Computation, v. 13, 2001.

Tan, P.N.; Steinbach, M.; Kumar, V, Introdução ao Data Mining: Mineração De Dados. Editora Ciência Moderna, 2009.

Vapnik, V. N. The nature of statistical learning theory. Springer-Verlag New York, Inc, 1995.

Vargas, R.; Munaro, C.; Ciarelli, P. Um Método para Detecção de Causalidade de Granger com Seleção de Regressores. XXI Congresso Brasileiro de Automática. Vitória: Sociedade Brasileira de Automática. 2016. p. 3397-3402.

Vasconcelos, I. O., Detecção móvel e online de anomalia em múltiplos fluxos de dados: Uma abordagem baseada em processamento de eventos complexos para detecção de comportamento de condução, Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática, 2017.