Fault detection in a ball mill circuit using principal component analysis

Pedro Henrique Trindade Dias Cabral¹ ^[1] Leonardo Junior Fernandes Campos^{1*} ^[1] André Luiz Alvarenga Santos¹ ^[1] Douglas Batista Mazzinghy¹ ^[1]

Abstract

Quality engineering is a fundamental aspect of production systems that require minimal variability. However, the large number of variables analyzed in industrial processes and the correlation between them requires a more robust technique such as Principal Component Analysis (PCA) for process monitoring. The aim of this work was to develop a fault detection system, using the R programming language, based on PCA in order to improve the quality of products in a milling circuit and facilitate decision-making. The developed algorithm was validated using the Benchmark Tenessee case study, whose maximum deviations were less than 5%. Applying this algorithm to a real case study made it possible to detect a pulp box overflow fault, with very low T^2 (0.012) and Q (0.026) values for the non-detection rate (TND).

Keywords: Grinding circuit; Multivariate statistic; Process control; Principal Component Analysis.

1 Introduction

Quality has been a crucial factor for the processing industry around the world over the decades. Given the competition, globalization, and technological aspects, to name a few, its role has increased exponentially in recent times. More specifically, quality engineering is the set of operational, managerial, and engineering activities that a company uses to ensure the quality characteristics of a product, with minimal variability. As variability can only be described in statistical terms, statistical methods play a central role in quality control, which is based on specification limits. In this sense, Statistical Process Control (SPC) emerged as a powerful set of problem-solving tools in the last century [1,2].

The control chart is one of the main SPC techniques. The most representative is the Shewhart control chart, which has been widely used to monitor quality characteristics in manufacturing. System monitoring tasks have three main steps. Namely, the detection of an undesired condition, called fault, diagnosis of this condition, and interventions to recover the process [3]. A disadvantage of the Shewhart control chart is its univariate nature. Thus, the correlation between the variables, an inherent characteristic of continuous industrial processes, is not considered. Coupled with the high number of variables measured mostly nowadays, principal component analysis (PCA) is therefore commonly applied for process monitoring [4-6]. Principal component analysis (PCA) belongs to the area of multivariate statistics, whose principle is to reduce the dimensionality of a problem by keeping most of its variance [7]. The greater the correlation between the variables, the greater the dimensionality reduction. Thus, it can deal with highly correlated multivariate systems, as is the case in the process industry. This characteristic explains the large number of applications in process monitoring. PCA can be seen as a multivariate extension of the univariate Shewhart control chart. The following works are some examples in the mining sector [8-13].

The objective of this paper was to develop a PCA-based fault detection system aimed at improving quality and also decision-support tool in a grinding circuit from Brazilian iron ore mine. This paper was presented as a case studies: the first one was used to develop the script and its validation, while the second one presents the application of the script validated and developed on the first case.

2 Principal component analysis (PCA)

A brief description of the principal component analysis (PCA) formulation is described below using a matrix representation. More information can be found in fundamental literature [2,7].

*Corresponding author: leocampos@demin.ufmg.br

Addresses: dias.phtc@gmail.com; asantos@demin.ufmg.br; dmazzinghy@demin.ufmg.br



^{2176-1523 © 2025.} Cabral et al. Published by ABM. This is an Open Access article distributed under the terms of the Creative Commons Attribution license (https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

¹Departamento de Engenharia de Minas, Escola de Engenharia, Universidade Federal de Minas Gerais, Belo Horizonte, MG, Brasil.

Cabral et al.

Consider a matrix $X_{[n,p]}$ with n records and p process variables. The records are related to the values assumed by the variables in the sampling interval and the variables can be, for example, flow, pressure, temperature, concentration, etc. One of the first actions for this database is to normalize each variable to zero mean and unit variance. This is important because the magnitudes of the variables can span a considerable range. Thus, the normalized matrix $Z_{[n,p]}$ is used instead of the matrix $X_{[n,p]}$ containing the actual values of the variables. In this normalization process, the correlation matrix $R_{[n,p]}$ of Z is obtained. Using the characteristic polynomial equation, the correlation matrix can be described by a set of pairs of eigenvalues and pairs of eigenvectors. The eigenvalues are given by the vector $\Lambda_{[p,1]}$ and the eigenvectors by the matrix $W_{[p,p]}$, where the eigenvectors are organized in columns ($w_{[p,1]}$). Both are arranged in descending order according to the values in Λ .

From a practical point of view, the eigenvectors are the axes of the new coordinate system, derived from the orthogonal rotation of the original coordinate system defined by X. They are called principal components (CPs). They are linear combinations of the original variables but not correlated with each other.

The corresponding eigenvalues $(\lambda_i, i=1,2,...,p)$ are the variances of each of the principal components. The original total variance is also preserved in the new coordinate system, as shown by Equation 1, where σ_i^2 is the variance of the ith-variable. The eigenvectors are defined in such a way as to maximize the original variance explained by them. In other words, the first principal component explains the maximum possible variance, the second, the remaining possible variation, and so on. Thus, few components can describe the main characteristics of the process under analysis.

This feature opens an opportunity to reduce its dimension by using only a few components (k) instead of the total number of original variables (p), with $k \ll p$. The number of components to retain in the model is generally defined by setting the desired amount of variation to be explained. The remaining components are often associated with process noise.

$$\sum_{i=1}^{p} \lambda_i = \sum_{i=1}^{p} \sigma_i^2 \tag{1}$$

The sample points in the new coordinate system are called scores $(S_{[n,p]})$. They are calculated according to Equation 2. It is noteworthy that the spatial arrangement of the original data is preserved in this new coordinate system. The original normalized values $(Z_{[n,p]})$ can be recovered as shown in Equation 3.

$$S_{[n,p]} = Z_{[n,p]} \cdot W_{[n,p]}$$
⁽²⁾

$$Z_{[n,p]} = S_{[n,p]} \cdot W_{[p,p]}^{-1}$$
(3)

Tecnol Metal Mater Min., São Paulo, 2025;22:e3214

Two metrics are commonly employed in PCA-based process monitoring applications. The first is given by the T^2 statistic, shown in Equation 4. This metric uses the scores (s_i, i=1,2, ..., k) of only the first k components, thus describing the main characteristics of the process under analysis. The T^2 value is calculated for each sample so that it can be tracked over time in a monitoring task. Relatively larger deviations are an indication that the process is moving away from normal operation.

$$\Gamma^{2} = \sum_{i=1}^{k} \frac{s_{i}^{2}}{\ddot{e}_{i}}$$
(4)

The second metric is given by the Q statistic. Its calculation is analogous but using the remaining (p-k) components, as shown in Equation 5. The Q value is also calculated for each sample and, in addition to the information provided by T², relatively larger changes over time usually mean changes in the spatial noise correlation structure, as mentioned earlier.

$$Q = \sum_{i=k+1}^{p} s_i^2 \tag{5}$$

As the T^2 and Q statistics are multivariate in nature, they lead to multivariate process monitoring. For the definition of its multivariate control charts, control limits must be established. These limits can be calculated using the wellknown significance level (α), but the usual non-normality of process data can be a burden for this. A commonly used workaround is the percentile approach. For example, given a false alarm rate (FAR) of 1% [3], upper control limits are given by the 99th percentile using normal data. The lower control limits for both statistics are zero.

3 Case studies

3.1 Tennessee benchmark problem: validation of script

The Tennessee Benchmark has often been used in process engineering to develop fault detection systems. In this work, the Tennessee Benchmark was used to validate the algorithm in R related to the calculations used in the PCA and compared to the results presented in the literature by [3].

This Benchmark Tennessee case study is based on a real industrial process, the aim of which is to obtain two liquid products from four gaseous reactants. The process takes place through a set of four irreversible and exothermic chemical reactions that have an inert component and a by-product.

The Benchmark Tennessee case study system consists of five main pieces of equipment: a reactor, a condenser, a vapor-liquid separator, a recycle compressor and a stripper column. The system is structured in such a way as to allow the acquisition of 53 measurements which include process variables, manipulated variables and laboratory parameters, and are collected at a sampling interval of three minutes. The Tennessee Benchmark has 22 data sets, of which 21 are data from faulty operations and 1 data set corresponds to the normal operating condition. For practical purposes, Gaussian noise is introduced in all measurements. More detailed information on the Tennessee Benchmark can be found in [14].

3.2 Grinding circuit: application of script

The data set used in this case study refers to the operation of a grinding circuit in an iron ore beneficiation plant of a mining company in Brazil. Ore beneficiation consists of a set of operations to reduce particle size and concentrate the mineral particles of interest. The concentration operations remove the (unwanted) gangue minerals from the ores, concentrating the minerals of interest. Figure 1 shows a closed grinding circuit with a mill and hydrocyclone classifier typical of grinding operations, as well as the main system variables.

The grinding system consists of a ball mill (power: 2,486 hp; dimensions: 6.1 m x 4.4 m, and speed: 15.18 rpm), a slurry box and a classification system made up of high-frequency screens (dimensions: 1.22 m x 3.04 m; polyurethane screen opening: 0.150 mm; installed power: 2.5 hp per set of 5 screens). In the system under study, particle size distribution and solids concentration are the two quality parameters, and both affect the feed to the grinding circuit.

Table 1 contains the set of eight operational variables considered in the case study, which were selected on the basis of information from the literature and the empirical knowledge of the team at the mineral processing plant. The variables Y1-MV and Y2-MV are writing variables, and the others are process

and disturbance variables. The variables X1-PV and X2-PV are the ones of greatest interest, as they are indicators of the quality of the ore processing process. The data set corresponds to 19 days of continuous circuit operation, sampled every 5 seconds, totalling 259189 records per variable.

4 Metodology

Figure 2 shows the three main steps of the methodology. Each one is described below. This methodology was applied to both case studies. After its validation with the Tennessee benchmark problem (case study 1), it was applied to the grinding circuit (case study 2).

 Model identification: This step is responsible for obtaining the PCA model, characteristic of normal operating conditions. A parameter to be defined in this phase is the number (k) of principal components to be retained in the model.

Table 1. Variables of the grinding circuit of the case study

Variable	Code	Unit
Circuit feedrate	Y1-MV	t/h
Water flow rate added to slurry box	Y2-MV	m³/h
Slurry solids concentration	X1-PV	%
Particle size at 100#	X2-PV	#
Slurry box level	X3-DV	%
Slurry density	X4-DV	t/m ³
Circulating load	X5-DV	%
Mill efficiency	X6-DV	%



Figure 1. Grinding circuit flowsheet of the case study.



Figure 2. Methodology steps.

- Control limits definition: This step calculates an upper control limit (UCL) for each monitoring statistic, namely, T² and Q. For this, it is necessary to define the desired false alarm rate (FAR), which was equal to 1% in this work as adopted by Russell et al. [3]. Considering the usual non-normality in process data, the use of the well-known significance level (α) can lead to unsatisfactory results. A common solution in this case is the percentile approach; in this case, the 99th percentile. This calculation is based on the T² and *Q* values calculated from the previous PCA model using fault-free data. Thus, at this point, there is a multivariate control chart for each monitoring statistic.
- Fault detection: The fault data was then fed into the previous PCA model and the resulting T² and Q values were plotted against the respective control charts. For evaluation purposes, Missed Detection Rate (MDR; in %) was used as a performance metric. The values of T² and Q are expected to exceed the upper control limit after the fault occurs. Detection is then considered missed when it does not occur. The MDR calculation is shown in Equation 6, where the numerator refers to the number of records below the control limit after the fault occurs, and the denominator is given by the total number of records in the time interval under fault. The closer the MDR to zero, the better.

$$MDR = \frac{\text{Number of missed register}}{\text{Total number or registers under fault}}$$
(6)

5 Results and discussions

5.1 Tennessee benchmarck problem: validation of script

Table 2 shows the MDR results obtained from the PCA model for the T² and Q monitoring statistics. They also show the differences $(MDR_{This work} - MDR_{[3]})$ in relation to the work by [3]. Most of the difference values are below 1%,

as is desirable. The greater differences, in absolute terms, are equal to 4.1 (fault 21) and 2.8% (fault 5) for T^2 and Q, respectively, that is, below 5%. These results satisfactorily validate the PCA scripts written using the R programming environment for this work. This type of result is important before dealing with real-world problems, as is the second case study of this work.

Furthermore, it can be noted that for some cases both metrics consistently indicate fault (MDR values close to zero) (for example, fault 1), while for others only one metric is successful (for example, fault 4). In other cases, neither of the two metrics can identify the fault (MDR values close to one) (for example, fault 3). These results illustrate the greater difficulty in detecting some faults and that no system can recognize all types of faults in a process. The fault detection task is still a challenge for the process industry.

5.2 Grinding circuit: application of script

After validating the PCA scripts with the previous benchmark case study, the methodology was applied to a real case study of the grinding circuit of an iron ore beneficiation plant of a mining company in Brazil.

Table 3 shows the result of the PCA modelling for the grinding circuit using normal operating data. As the total number of original variables is relatively low, a numerical analysis of their weights (loads) in the model can be done more directly. It can be observed that the variables Y2-MV, X1-PV, X3-DV, X4-DV, and X5-DV have considerable weights (highlighted in bold) in the first principal component (CP1), which explain 36.3% of the variance of the process data. For the second principal component (CP2), which accounts for 22.2% of the explained variance, this occurs for Y1-MV and X7-DV. In the third principal component (CP3), which explains 12.0% of the original variance, X2-PV and Y2-MV again have relatively higher weights.

Thus, in a way, all eight original variables (Table 1) are covered by only three principal components (k = 3), which can explain 70.5% of the variance of the original data. This value is quite satisfactory in industrial engineering. Thus, the final PCA model to be used as a fault detection system

Table 2. Fault detection	results (comparison	with	Russell	et al.	[3])
						1 - 17

		\mathbf{T}^2		Q			
Fault —	MDR This work	MDR (Russell et al. [3])	Difference (%)	MDR This work	MDR (Russell et al. [3])	Difference (%)	
1	0.009	0.008	0.1%	0.004	0.003	0.1%	
2	0.021	0.020	0.1%	0.015	0.014	0.1%	
3	0.999	0.998	0.1%	0.999	0.991	0.8%	
4	0.963	0.956	0.7%	0.036	0.038	-0.2%	
5	0.746	0.775	-2.9%	0.774	0.746	2.8%	
6	0.013	0.011	0.2%	0.001	0	0.1%	
7	0.079	0.085	-0.6%	0.001	0	0.1%	
8	0.034	0.034	0.0%	0.026	0.024	0.2%	
9	0.999	0.994	0.5%	0.980	0.981	-0.1%	
10	0.644	0.666	-2.2%	0.649	0.659	-1.0%	
11	0.788	0.794	-0.7%	0.344	0.356	-1.2%	
12	0.023	0.029	-0.7%	0.026	0.025	0.1%	
13	0.061	0.060	0.1%	0.046	0.045	0.1%	
14	0.121	0.158	-3.7%	0.001	0	0.1%	
15	0.980	0.988	-0.8%	0.971	0.973	-0.2%	
16	0.829	0.834	-0.5%	0.744	0.755	-1.1%	
17	0.254	0.259	-0.5%	0.100	0.108	-0.8%	
18	0.114	0.113	0.1%	0.101	0.101	0.0%	
19	0.999	0.996	0.3%	0.858	0.873	-1.6%	
20	0.703	0.701	0.2%	0.546	0.550	-0.4%	
21	0.695	0.736	-4.1%	0.576	0.570	0.6%	

Table 3. PCA model for the second case study

Variable	CP1	CP2	СР3	CP4	CP5	CP6	CP7	CP8
Y1-MV	0.0514	0.9035	-0.0587	-0.1466	0.0588	0.0602	0.3861	0.0016
Y2-MV	-0.5603	0.0434	0.6445	-0.3700	-0.0354	-0.3610	0.0174	-0.0003
X1-PV	-0,8748	0.2049	-0.2686	0.2874	-0.1138	-0.1358	-0.0164	-0.0799
X2-PV	0.2740	0.2319	0.5589	0.7317	-0.1418	0.0510	0.0246	0.0035
X3-DV	-0.6199	-0.1229	0.1482	0.1473	0.7323	0.1427	0.0188	0.0019
X4-DV	-0.8999	0.1331	-0.2723	0.2391	-0.1456	-0.1145	-0.0094	0.0819
X5-DV	-0.7410	-0.1024	0.2421	-0.2245	-0.3147	0.4820	0.0024	-0.0049
X6-DV	0.0679	0.9037	0.0329	-0.1292	0.0905	0.0658	-0.3853	0.0033
Variance (%)	36.3%	22.2%	12.0%	11.5%	8.8%	5.3%	3.7%	0.2%
Cumulative variance (%)	36.3%	58.5%	70.5%	82.0%	90.8%	96.1%	99.8%	100.0%

in the grinding circuit is composed of PC1, PC2 and PC3. Table 3 shows the MDR results obtained from the PCA model for the T^2 and Q monitoring.

Next, given the 1% false alarm rate (FAR) as before, the upper control limits (UCL) for the T^2 and Q statistics were calculated using the 99th percentile procedure. They are about 20 and 200, respectively.

Then, a fault detected on November 7, 2015 by the plant team responsible for operating the grinding circuit of the case study was selected to be used in this work. It started around 6:50 am and ended around 8:30 am. This event referred to an overflow in the pump box when analyzing the rapid increase of the tank level without increasing the pump frequency (Figure 3).



Figure 3. Pump box level and pump frequency data.



Figure 4. Multivariate control charts for (A) T² and (B) Q statistics, given the pump box fault in the grinding circuit.

The multivariate control charts for the T^2 and Q monitoring statistics corresponding to this fault period are shown in Figure 4. The MDR values for them were respectively equal to 0.012 and 0.026. These significantly low values are an indication that the fault was detected within a reasonable amount of time. The sooner a fault is detected, the greater the chance of recovering the process and mitigating potential losses. Furthermore, it can be seen that there were no false alarms, which is crucial for the reliability of a fault detection system.

6 Conclusions

The algorithm developed for fault determination proved to be effective during validation with the Benchmark Tennessee case study, showing fault detection deviations of less than 5% in relation to literature data. The application of this algorithm to industrial data from the iron ore grinding circuit showed that only three principal components are sufficient to explain approximately 70% of the variance in the data and, as a result, the size of the analysis was reduced from 8 to 3. Thus, the eight original variables analysed are covered by the three principal components. The evaluation of the T^2 and Q statistics for a known case of failure reported by operators made it possible to identify the start and end of the failure by extrapolating the upper limits set for the statistics.

Acknowledgements

The authors express their gratitude to all the collaborators of this study and CAPES for financially supporting this study.

References

- 1 Jackson JE. Quality control methods for several related variables. Technometrics. 1959;1(4):359-377. http://doi.org/ 10.1080/00401706.1959.10489868.
- 2 Montgomery DC, Farias AML. Introdução ao controle estatístico de qualidade. 7[°] ed. Rio de Janeiro: LTC; 2016.
- 3 Russell EL, Chiang LH, Braatz RD. Fault detection in industrial processes using canonical variate analysis and dynamic principal component analysis. Chemometrics and Intelligent Laboratory Systems. 2000;51(1):81-93. http:// doi.org/10.1016/S0169-7439(00)00058-7.
- 4 Severson K, Chaiwatanodom P, Braatz RD. Perspectives on process monitoring of industrial systems. Annual Reviews in Control. 2016;42:190-200. http://doi.org/10.1016/j.arcontrol.2016.09.001.
- 5 Yin S, Ding SX, Xie X, Luo H. A review on basic data-driven approaches for industrial process monitoring. IEEE Transactions on Industrial Electronics. 2014;61(11). http://doi.org/10.1109/TIE.2014.2301773.
- 6 Qin SJ. Statistical process monitoring: basics and beyond. Journal of Chemometrics. 2003;17:480-502. http://doi. org/10.1002/cem.800.
- 7 Mingoti SA. Análise de dados através de métodos de estatística multivariada: uma abordagem aplicada. Belo Horizonte: Editora UFMG; 2005.
- 8 Marais C, Aldrich C. Estimation of platinum flotation grades from froth image data. Minerals Engineering. 2011;24(5):433-441. http://doi.org/10.1016/j.mineng.2010.12.006.
- 9 Pumure I, Renton JJ, Smart RB. The interstitial location of selenium and arsenic in rocks associated with coal mining using ultrasound extractions and principal component analysis (PCA). Journal of Hazardous Materials. 2011;198:151-158. http://doi.org/10.1016/j.jhazmat.2011.10.032.
- 10 Pandey B, Agrawal M, Singh S. Assessment of air pollution around coal mining area: Emphasizing on spatial distributions, seasonal variations and heavy metals, using cluster and principal component analysis. Atmospheric Pollution Research. 2014;5(1):79-86. http://doi.org/10.5094/APR.2014.010.
- 11 Wang X, Liu J, Carranza EJM, Wang J, Wang G, Zhai D, et al. A combined approach using spatially-weighted principal components analysis and wavelet transformation for geochemical anomaly mapping in the Dashui ore-concentration district, Central China. Journal of Geochemical Exploration. 2019;197:228-237. http://doi. org/10.1016/j.gexplo.2018.12.008.
- 12 Sun X, Zhou Y, Yuan L, Li X, Shao H, Lu X. Integrated decision-making model for groundwater potential evaluation in mining areas using the cusp catastrophe model and principal component analysis. Journal of Hydrology. Regional Studies. 2021;37:100891. http://doi.org/10.1016/j.ejrh.2021.100891.
- 13 Seyedrahimi-Niaraq M, Mahdiyanfar H, Mokhtari AR. Integrating principal component analysis and U-statistics for mapping polluted areas in mining districts. Journal of Geochemical Exploration. 2022;234:106924. http://doi. org/10.1016/j.gexplo.2021.106924.
- 14 Downs JJ, Vogel EF. A plant-wide industrial process control problem. Computers & Chemical Engineering. 1993;17(3):245-255. http://doi.org/10.1016/0098-1354(93)80018-I.

Received: 22 Jan. 2025 Accepted: 22 Abr. 2025

Editor-in-charge: André Carlos Silva 💿