

A Self Adaptive Genetic Based Algorithm for the Identification and Elimination of Bad Data

A. A. Hossam-Eldin, E. N. Abdallah, and M. S. El-Nozahy

Abstract—The identification and elimination of bad measurements is one of the basic functions of a robust state estimator as bad data have the effect of corrupting the results of state estimation according to the popular weighted least squares method. However this is a difficult problem to handle especially when dealing with multiple errors from the interactive conforming type. In this paper, a self adaptive genetic based algorithm is proposed. The algorithm utilizes the results of the classical linearized normal residuals approach to tune the genetic operators thus instead of making a randomized search throughout the whole search space it is more likely to be a directed search thus the optimum solution is obtained at very early stages(maximum of 5 generations). The algorithm utilizes the accumulating databases of already computed cases to reduce the computational burden to minimum. Tests are conducted with reference to the standard IEEE test systems. Test results are very promising.

Keywords—Bad Data, Genetic Algorithms, Linearized Normal residuals, Observability, Power System State Estimation.

I. INTRODUCTION

POWER system state estimation is the process of calculating a reliable estimate of the power system state vector composed of bus voltages angles and magnitudes from telemetered measurements on the system. This estimate of the state vector provides the description of the system necessary for the operation and security monitoring. However those telemetered measurements may have errors and those errors may affect the accuracy of the estimated values, thus many efforts were devoted for the issue of bad data identification and elimination. The most common approach is the use of the Linearized Normal Residuals (LNRs) in accordance with a statistical criterion (chi square test) however it showed many drawbacks as the identification procedure often can not pinpoint a single bad measurement but instead identifies a group of measurements some of which is bad. In such cases the group must be eliminated to eliminate the bad measurement. Moreover successive elimination of a group of measurements with the highest LNR may turn the network unobservable. Also, successive suppression of measurements with the highest LNRs may leads to suppression of the correct

measurements instead of the bad ones as shown in the famous case study described in [1].

Bad measurements may also be identified by using a non-quadratic state estimator in solving the state estimation problem such as the weighted absolute sum of residuals [2] owing to the automatic bad data rejection property of these estimators [3]. However it is subject to convergence problems and constitutes a huge computational burden.

Reference [4] formulated the bad data identification problem as a combinatorial problem and solved it by the branch-and-bound method. However Branch-and-bound method may be replaced by faster non deterministic methods.

Reference [5] used artificial neural networks to solve the bad data identification combinatorial problem. Tabu search was also proposed in [1].

Genetic algorithms were introduced in [6] to solve the bad data identification combinatorial problem. The authors tried 3 different versions of genetic algorithms; the basic genetic algorithm, the micro-genetic algorithm and the single individual genetic algorithm. The number of state re-estimations was reduced by using a database of already computed cases and a filtering mechanism was applied to skip non promising solutions; however the computational burden was still too high as it requires 100 generations to reach the optimum solution which makes it unsuitable to be used online.

In this paper a new objective function for handling the bad data identification problem is proposed, methods to optimize it were implemented using a self adaptive genetic based algorithm. The algorithm utilizes the results of the classical linearized normal residuals approach to self tune the genetic operators thus instead of making a randomized search throughout the whole search space it is more likely to be a directed search thus the optimum solution is obtained at very early generations(maximum of 5 generation).

The concept of the accumulating databases introduced in [6] is also adopted here to reduce the number of state re-estimations to minimum .

Finally tests were conducted on standard IEEE test systems to test the robustness of the proposed algorithm and the results showed its high efficiency in identifying multiple errors from different types at very low computational cost.

II. PROBLEM FORMULATION

According to [4] the bad data identification problem is an optimization problem with a combinatorial nature.

Manuscript received August 28, 2008.

Authors are with the department of Electrical Engineering, Faculty of Engineering Alexandria University, P. O. Box: 21544, Alexandria, Egypt.

For a system with m measurements the suspect measurement set is represented by an m -dimensional decision vector, b , in which:

$$\begin{aligned} b_i &= 1 && \text{if the } i^{\text{th}} \text{ measurement is a bad data.} \\ b_i &= 0 && \text{if the } i^{\text{th}} \text{ measurement is good.} \end{aligned}$$

Thus creating 2^m possible decision vectors where each decision vector will represent a possible combination of good and bad measurements. After each decision vector set (b) is formed, the state estimation problem is solved using only the unexcluded measurements and the decision vector is assumed valid if the updated resulting sum of residuals $J[x^{\text{new}}(b)]$ is less than the updated chi square threshold $t_J(b)$ keeping into consideration that set (b) is sufficient to maintain the system observable. Those valid decision vectors are evaluated and the one with highest fitness is considered the optimum solution for the problem.

According to [5] the occurrence of bad data is rather unlikely, the valid decision vector that includes least number of bad data is the one with the highest fitness. The problem of identification of bad data was formulated as follows:

$$\text{Minimize } F(b) = \sum_{i=1}^m b_i \quad (1)$$

Subject to: Set (b) is observable
 $J[x^{\text{new}}(b)] < t_J(b)$

However this is contradicted by the fact that in the case of multiple bad data, it is not necessary to eliminate all the bad data present among the measurements set for the resulting decision vector to pass the chi square test, moreover as shown in the test case presented in [1] upon the insertion of two interacting bad data from the conforming type in the measurements set, the LNR approach identified a good measurement instead of the two bad ones and eliminated it, however the updated decision vector (with two bad data present and without the excluded good datum) passed the chi test.

Also considering that the lowest sum of residuals $[J(x)]$ is obtained when the measurements set is free from bad data the authors proposed a modified objective function as follows:

Minimize

$$F(b) = \sum_{i=1}^m b_i + [A_2 \times J(x^{\text{new}})] \quad (2)$$

Subject to:

$$\begin{aligned} \text{Set (d) is observable} \\ J[x^{\text{new}}(d)] < t_J(d) \end{aligned}$$

In our proposed form we aim not only to minimize the number of identified bad measurements among a certain measurements set that has to be excluded to pass the chi square statistical criterion, but we also aim to minimize the sum of residuals obtained by solving the state estimation problem using the rest of the unexcluded measurements keeping in mind not to exclude measurements that will render the system unobservable. A_2 is the quality fitness weight; it

assigns the weight given to the minimization of the measurement residuals on the expense of increasing the number of detected bad measurements. A small value for A_2 like 0.1 was found to give satisfactory results for both IEEE 6 bus bars and 14 bus bars test systems. i.e. if we have 2 possible decision vectors, set₁ (b) has 6 detected bad measurements and set₂ (b) has 5 detected bad measurements, both sets satisfied the specified constraints, set₂ (b) is always favored to set₁ (b) except if set₁ (b) results in reducing a measurements residual less than that introduced by set₂ (b) by at least $\frac{1}{A_2}$. However for other test systems the optimum value of A_2 could change and can be found by trial.

III. IMPLEMENTATION

A. Fitness Evaluation

The fitness function used in this paper is

$$\text{Fit}(x) = \begin{cases} F_{\max} - F(x) & , \text{ if } F(x) < F_{\max} \\ 1 & , \text{ otherwise} \end{cases} \quad (3)$$

Where

$F(x)$ is the raw fitness function = $-F(b)$

F_{\max} is the largest value of $F(x)$ in the current population .

The fitness function $\text{Fit}(x) = \frac{1}{F(b)}$ was adopted in [6] and

was tested but it showed lower computational efficiency than the chosen form.

B. Representation of Variables

Since each gene in the decision vector has only 2 states, either ($b_i = 1$ if the i^{th} measurement is in gross error or $b_i = 0$ if the i^{th} measurement is an accurate one) so the most appropriate representation of variables is the binary encoding.

C. Using the Results of the Classical Approach to Generate the Probabilistic Linearized Normal Residuals (PLNRs)

It is clear that as the LNR for a certain measurement increases its probability to be a bad measurement also increases, thus we introduced the concept of the Probabilistic Linearized Normal Residual (PLNR) as follows:

$$\text{PLNR}_i = \frac{|LNR_i|}{\sum_{j=1}^{N_m} |LNR_j|} \quad (4)$$

The PLNR for a certain measurement is in the range between 0 and 1 and it shows how probable a certain measurement is bad with respect to all other measurements. PLNRs will be used in the next sections to tune the genetic operators as well as generating the initial population.

D. Generating the Initial Population

Initial population is generated in a deterministic manner associated with a random part. Each chromosome (representing a decision vector) has number of genes equals to the number of measurements where each gene has 2 possible as stated before.

We developed the following formula for the generation of initial population:

$$b_i = \begin{cases} 1 & \text{if } PLNR_i > \text{mean}(PLNRs) + [\text{std}(PLNRs) * A_1 * \text{rand}(1)] \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

Where $\text{mean}(PLNRs)$ is the mean value for all the $PLNRs$ in the decision vector, $\text{std}(PLNR)$ is their standard deviation, $\text{rand}(1)$ is a random number between 0 and 1 to guarantee the diversity required for the genetic operation.

A_1 is the window factor; it determines the range of suspect of the presence of bad measurements. At the beginning the window factor is set to a high value (15 for example) thus we suspect only measurements whose $PLNRs$ are greater than the mean value of all $PLNRs$ by 15 times the standard deviation to be bad. After the generation of the first chromosome (representing a possible decision vector) we check the validity of such chromosome (whether it satisfies the constraints or not); if not we generate another chromosome and check its validity. After 10 unsuccessful trials to generate a valid chromosome at the initial window ratio specified ($A_1=15$) it is revealant that we have a good reason to suspect the presence of more bad measurements than the ones included in such window so we have to increase our window of search (decrease A_1 to 14) and so on. This guarantees that we narrow the search space as much as possible to obtain the solution in early generations.

E. Handling Constraints

Since genetic algorithms (GAs) are directly applicable only to unconstrained optimization, it is necessary to use some additional methods that will keep solutions in the feasible region. During the past few years, several methods were proposed for handling constraints by GAs. Most of these methods are problem-dependent (i.e. specific algorithm has to be designed for each particular problem). The most popular approaches in GA community to handle constraints are the following techniques

- Penalty technique
- Rejection Technique
- Repairing technique

The technique adopted for each constraint introduced in Eq 2 depends on the constraint itself:

Observability constraint: two techniques are adopted here; if the system is unobservable due to the removal of 1 critical measurement only, the resulting decision vector can be repaired by finding out the critical measurement identified as a bad one and assigning it to be valid measurement again i.e. changing $b_i = 1$ into $b_i = 0$ (repairing technique). This is because critical measurements can be identified easily as they are the ones with zero measurement residuals.

However if the system is unobservable due to insufficiency of measurement devices or simultaneous removal of a critical pair or critical k-tuple measurements , repairing techniques involves finding out the unobservable islands and the locations measurement devices are to be placed to restore observability. This can be carried out by numerical observability methods but it will introduce great computational burden, so instead the simple observability algorithm based on the topological observability concept described in [7] will be adopted here to check whether the system is observable for a certain decision vector or not and if the resulted decision vector was found to be insufficient to make the system fully observable it will be rejected and a new one is generated which will save a huge computational burden.

Chi test constraint: the resulted decision vector is checked to find out whether it passes the chi test i.e. the updated sum of residuals after solving the state estimation problem using only the good measurements $J[x^{\text{new}}(b)]$ is less than the updated chi square threshold $t_j(b)$, if not it will be rejected and another decision vector will be generated instead. This insures that all chromosomes existing in different populations are valid ones, which resulted in reducing the computational burden significantly and the optimum solution is obtained at very early generations.

F. Genetic Operators

i) *Population size:* The use of the rejection technique for decision vectors that violate the problem constraints allowed us to use a reduced population size. The population size adopted here is equal to the square root of the number of measurements.

ii) *Maximum generations:* One of the greatest advantages of the proposed algorithm is that the solution is obtained at early generations, so a maximum number of generations equals to 5 is enough to guarantee obtaining optimum solution even for large power systems. In many of the test cases the optimum solution was obtained before the 5th generation as will be shown.

iii) *Crossover probability:* Since we used reduced population size, it is recommended to increase the crossover probability to favor exchanging genetic properties between different chromosomes in the same generation, that's why a fixed crossover probability of unity was used.

iv) *Selection:* Roulette wheel selection is adopted here.

v) *Mutation probability:* Mutation is responsible of producing random changes to the values of different genes in a chromosome.

A variable mutation probability is proposed as follows:

$$b_i = \begin{cases} 1 & , \text{if } b_i = 0 \& \text{rand}(1) < PLNR_i \\ 0 & , \text{if } b_i = 1 \& \text{rand}(1) > PLNR_i \end{cases} \quad (6)$$

This ensures that each gene is more probably to go mutation (its value is more probable to change from 0 to 1) as its $PLNR$ increases (as it is more probable to be a bad measurement) and vice versa .The use of the variable mutation probability

introduced above caused great improvement in the performance of the suggested algorithm.

G. Computational Burden Reduction Techniques

It is clear that running the state estimation subroutine for a certain decision vector constitutes a major load of the computational burden, so the concept of accumulating database was used in this algorithm to reduce the computational burden.

Two databases are created, one that includes all the evaluated decision vectors that satisfies all the system constraints and the other includes the invalid ones. Before running the state estimation subroutine for a newly generated decision vector it is compared first with those stored in the 2 databases; if found there is no need to re-estimate the system states, if not found in any of them it will be evaluated and added to the appropriate database.

Also using the rejection technique for decision vectors that violate the observability constraint saves the time and computational load that would have been consumed in determining the unobservable islands and finding the necessary measurements among the assumed bad ones that have to be added to restore observability.

IV. TESTS AND RESULTS

The algorithm was implemented in a software package using the MATLAB Programming language. It was tested on some standard IEEE test systems. All tests were carried out on a 3.6 GHz Pentium IV, 2 Giga Bytes Ram PC.

A. Tests on Standard IEEE 6 Bus Bars Test System

TABLE I

TESTS 1 & 2 ON THE STANDARD IEEE 6 BUS BARS TEST SYSTEM

Test Number		Test 1	Test 2
Bad Measurements	True	$P_1 = 1.079$ $Q_1 = 0.16$	$P_2 = 0.5$ $P_{21} = -0.278$
	Inserted	$P_1 = -1.131$ $Q_1 = -0.202$	$P_2 = 1.3$ $P_{21} = -0.75$
	Estimated	$P_1 = 1.1043$ $Q_1 = 0.1698$	$P_2 = 0.474$ $P_{21} = -0.286$
No. of detected bad data		2	2
Evaluated points	Valid	52	58
	Invalid	17	89
Data Hits	Valid	76	128
	Invalid	28	513
Execution time(seconds)		3.562	3.66
Gen. at which optimum solution is obtained		1	1
Type of bad data inserted		Multiple interacting non conforming	Multiple interacting non conforming
Updated sum of residuals $J(x^{new})$		40.6470	39.6416

TABLE II
TESTS 3 & 4 ON THE STANDARD IEEE 6 BUS BARS TEST SYSTEM

Test Number	Test 3	Test 4
Bad Measurements	True	$P_{14} = 0.436$ $Q_{14} = 0.201$ $P_{25} = 0.155$ $Q_{25} = 0.154$ $P_{36} = 0.438$ $Q_{36} = 0.607$ $P_{53} = -0.18$ $Q_{53} = -0.261$
	Inserted	$P_{14} = -0.389$ $Q_{14} = -0.212$ $P_{25} = -0.174$ $Q_{25} = -0.22$ $P_{36} = -0.433$ $Q_{36} = -0.583$ $P_{53} = 0.251$ $Q_{53} = 0.299$
Estimated		$P_2 = 0.5$ $P_{24} = 0.331$ $Q_2 = 0.744$ $Q_{24} = 0.461$
		$P_2 = 0.968$ $P_{24} = 0.656$ $Q_2 = 1.44$ $Q_{24} = 0.766$
No. of detected bad data	8	4
Evaluated points	Valid	70
	Invalid	294
Data Hits	Valid	205
	Invalid	176
Execution time(seconds)	4.05	3.78
Gen. at which optimum solution is obtained	3	5
Type of bad data inserted	Multiple interacting non conforming	Multiple interacting conforming
Updated sum of residuals $J(x^{new})$	33.2812	38.5147

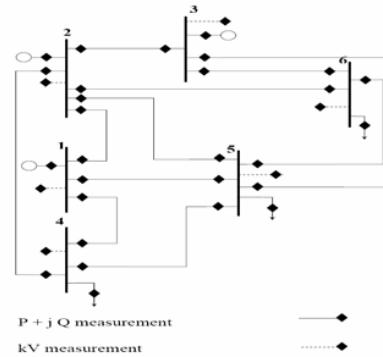


Fig. 1 IEEE 6 Bus bars system test

B. Tests on Standard IEEE 14 Bus Bars Test System

TABLE III

TESTS 1 & 2 ON THE STANDARD IEEE 14 BUS BARS TEST SYSTEM

Test Number		Test 1	Test 2
Bad Measurements	True	$P_1=2.3227$ $Q_1=-0.225$ $P_{42}=-0.5448$ $P_{45}=-0.6108$	$V_2=1.045$ $V_3=1.01$ $V_4=1.019$
	Inserted	$P_1=0.2$ $Q_1=-2$ $P_{42}=-1.2$ $P_{45}=-2.3$	$V_2=0.5$ $V_3=1.4$ $V_4=0.8$
	Estimated	$P_1=2.259$ $Q_1=-0.2343$ $P_{42}=-0.539$ $P_{45}=-0.6025$	$V_2=1.0469$ $V_3=1.0114$ $V_4=1.0209$
No. of detected bad data		4	3
Evaluated points	Valid	97	87
	Invalid	461	16
Data Hits	Valid	44	149
	Invalid	709	80
Execution time(seconds)		51.0875	11
Gen. at which optimum solution is obtained		1	2
Type of bad data inserted		Multiple interacting non conforming	Multiple non interacting
Updated sum of residuals $J(x^{new})$		3.0877	3.2682

TABLE IV

TESTS 3 & 4 ON THE STANDARD IEEE 14 BUS BARS TEST SYSTEM

Test Number		Test 3	Test 4
Bad Measurements	True	$P_2=0.1847$ $P_{24}=0.5616$ $P_{25}=0.4153$	$P_2=0.1847$ $P_{24}=-1.5242$
	Inserted	$P_2=2$ $P_{24}=1.8$ $P_{25}=1.5$	$P_2=1.2$ $P_{21}=-3.5$
	Estimated	$P_2=0.2706$ $P_{24}=0.567$ $P_{25}=0.423$	$P_2=0.2692$ $P_{21}=-1.4617$
No. of detected bad data		3	2
Evaluated points	Valid	90	88
	Invalid	346	48
Data Hits	Valid	142	260
	Invalid	699	66
Execution time(seconds)		46.469	13.563
Gen. at which optimum solution is obtained		1	1
Type of bad data inserted		Multiple interacting conforming	Multiple interacting conforming
Updated sum of residuals $J(x^{new})$		2.2498	2.3708

TABLE V
TESTS 5 & 6 ON THE STANDARD IEEE 14 BUS BARS TEST SYSTEM

Test Number	Test 5	Test 6
Bad Measurements	True	$P_{15}=0.7556$ $Q_{15}=0.0087$ $P_{24}=0.5616$ $Q_{24}=-0.042$ $P_{25}=0.4153$ $Q_{25}=-0.010$ $P_{42}=-0.5448$ $Q_{42}=0.0203$
	Inserted	$P_{15}=2.3$ $Q_{15}=1.9$ $P_{24}=3.7$ $Q_{24}=-2.4$ $P_{25}=1.8$ $Q_{25}=-2.2$ $P_{42}=-3.7$ $Q_{42}=2.1$
	Estimated	$P_{15}=0.7335$ $Q_{15}=0.0066$ $P_{24}=0.55$ $Q_{24}=-0.0399$ $P_{25}=0.4066$ $Q_{25}=-0.0081$ $P_{42}=-0.5339$ $Q_{42}=0.0164$
No. of detected bad data	8	6
Evaluated points	Valid	111
	Invalid	639
Data Hits	Valid	272
	Invalid	194
Execution time(seconds)	56.281	47.328
Gen. at which optimum solution is obtained	4	4
Type of bad data inserted	Multiple interacting non conforming	Multiple interacting conforming
Updated sum of residuals $J(x^{new})$	2.5277	3.1389

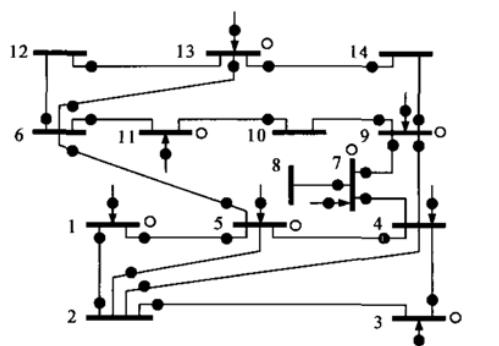


Fig. 2 IEEE 14 Bus bars test system

V. DISCUSSION

The proposed algorithm showed supreme efficiency in identifying all the bad measurements present in a certain measurements set owing to the robustness of the modified objective function presented in equation 2 even when introducing eight bad measurements (as in tests 3 and 5).

The algorithm was able to identify only the bad measurements and not a group of measurements some of which is bad which is the case in the classical approaches.

The enhancements in the genetic operation made it possible to obtain the optimum solution at very early stages of the program, in many cases solution was obtained in the first generation which proves the efficiency of the proposed technique used to generate the initial population.

Moreover the use of the accumulative databases concept reduced the execution time of the program that it becomes much faster than the algorithms described in the literature.

We notice that the sum of residuals ($J(x^{new})$) and consequently the errors in the estimated state variables obtained for tests applied on 14 bus bars test system are much less than those obtained for the 6 bus bars test system. This can be justified that the degree of freedom in 14 bus bars test system is 95 (122 available measurements are used to estimate 27 state variables) while for the 6 bus bars test system the degree of freedom is 51 (62 measurements are used to estimate 11 state variables) and thus for the 14 bus bars test system we have more redundancy, similarly we conclude that for practical systems (118 and 300 bus bars test systems) we will have more redundancy and thus even better results can be obtained.

Power systems tend to be static in nature (its operating conditions do not change rapidly) thus the run time of the program (less than 1 minute) makes it possible to use our proposed algorithm in practical Energy Management Systems (EMS) in the online mode to identify bad data especially when using a supercomputer instead of a PC.

VI. CONCLUSION

A new robust algorithm for handling the bad data identification problem was introduced in this paper. A modified objective function is proposed where we aim not only to minimize the number of identified bad measurements among a certain measurements set that has to be excluded to pass the chi square statistical criterion, but we also aim to minimize the sum of residuals obtained by solving the state estimation problem using the rest of the unexcluded measurements keeping in mind not to exclude measurements that will render the system unobservable. We proposed a self adaptive genetic based technique that uses the results obtained from the chi square statistical criterion to calibrate the genetic algorithm (GA) parameters in order to pinpoint suspected measurements.

The algorithm was tested on the IEEE standard 6 bus bars test system and 14 bus bars test system and it showed supreme performance.

The self adaptation methods and the accumulative databases concept introduced in the proposed algorithm resulted in a significant reduction in the computational time as the

algorithm performs a directed search rather than a randomized search throughout the whole solution space and thus convergence can be attained at early generations (maximum of 5 generations) using a small population size (square root the number of measurements). Moreover it showed 100% efficiency in eliminating all the bad data present in a certain measurement set.

REFERENCES

- [1] A. Monticelli, *State estimation in electric power systems. A generalized Approach*, Boston: Kluver Academic Publishers, 1999, ch.9.
- [2] A. A. Abur, "A bad data identification method for linear programming state estimation", *IEEE Trans. Power Systems*, vol. PWRS-5, No. 3, pp.894-901, August 1990
- [3] W. W. Kotiuga and M. Vidyassagar, "Bad data rejection properties of weighted least absolute value techniques applied to static state estimation", *IEEE Trans. Power Apparatus and Systems*, vol. PAS-101, No.2, pp. 511-523, May 1991.
- [4] A. Monticelli, F. Wu and M. Yen, "Multiple bad data identification for state estimation by combinatorial optimization", *IEEE Trans. Power Delivery*, vol. PWRD-1, No.3, pp. 361-369, July 1986.
- [5] N. H. Abbasy and W. El-Hassawy, "Power system state estimation: ANN application to bad data detection and identification", *Proc. 4th IEEE AFRICON Conference*, September 1996, vol.2, pp.611-615.
- [6] S. Gastoni, G. P. Granelli and M. Montagna, "Multiple bad data processing by genetic algorithms", *Proc. IEEE Bologna PowerTech Conference*, June 2003.
- [7] G. R. Krumpholz, K.A. Clements and P.W. Davis, "Power Systems Observability: A Practical Algorithm Using Network Topology", *IEEE Transactions on Power Apparatus and Systems*, vol. PAS-99, No.4, August 1980, pp. 1534-1542.