

Predictive Quality Control^{☆☆}

C.V. Radhakrishnan^{a,1,**}, K. Bazargan^{a,b,2}

^a*River Valley Technologies, SJP Building, Cotton Hills, Trivandrum, Kerala, India 695014*

^b*River Valley Technologies, 9, Browns Court, Kennford, Exeter, United Kingdom*

Abstract

Abstract, should normally be not longer than 200 words.

Keywords: quadruple exciton, polariton, WGM

PACS: 71.35.-y, 71.35.Lk, 71.36.+c

2008 MSC: 23-557

1. Introduction

1.1. How is quality measured in present time?

Present day quality measurements are done using various metrological tools which quantify various shapes and/or form features of the part. These features are given their individual limits by the designer which help a metrologist in determining whether a part that has just been manufactured is fit to be deemed a good part or not. However, most of these tools are completely external equipments which need to be explicitly used over the part(s) to see the results. This usually calls for extra man-power and reduced throughput. So what manufacturers end up doing is that they keep an inspection frequency for non critical parts where quality checks on these parts are made at regular intervals until one sees an anomaly. How this anomaly is responded to again depends on what feature we are measuring.

The main problem with such a way of measurement is that :-

1. We have to make a compromise between the throughput and inspection interval.
2. As the quality measurement is isolated from part manufacturing, you can't detect a bad part unless it has already been manufactured.

[☆]This document is a collaborative effort.

^{*}Corresponding author

^{**}Principal corresponding author

Email address: cvr@river-valley.com (C.V. Radhakrishnan)

URL: <http://www.elsevier.com> (K. Bazargan)

¹This is the specimen author footnote.

²Another author footnote, but a little more longer.

1.2. How do we propose to measure it?

The main gist of the ideology we have applied is that, when a part is being made, the information of each cut made by the machine is reflected in some real-time aspect of the machine. It could be noise, power, vibration, loads, etc. So what we do is, we record these variables (in our case, power) real-time using appropriate tools and find a correlation between these readings and the quality of the part manufactured. This correlation can be tightened depending on the kind of accuracy one is looking for. It also has an advantage over the conventional approach in certain key aspects such as:-

1. The computational time is significantly less once a correlation model has been made. Also, buying computational time is cheaper than wasting producing time.
2. The feedback we get is near to real time and we could infact get early warnings as the part is being made.
3. The kind of investment needed for this is much lesser compared to buying precise instruments which, other than being pretty expensive, usually have only a single functionality.

2. Methodology*2.1. Getting Quality data*

As we need to correlate between power readings during a producing cycle and quality, we are essentially trying to create a model or create a procedure to generate a model which would take the power readings as input and give the quality as output. Thus, we needed some input quality data for a training set.

For this, we monitored the Ra values of a part, which underwent turning operation, at regular intervals (once in 25 parts for first 500 cycles and once in every 10 parts from there on) for 1030 parts. These measurements were made at two different faces and three different locations in each face. We observed a downward trend in the Ra values for the first 300 parts after which they were more or less constant (no particular trend) till part 3000 after which they again rose up and went outside the specification limit.

So we decided to use the trend evident in the first 300 parts for our model creation and verification. We created a 2 degree polynomial fit for the first three hundred parts using 12 readings as out control points.

Now this sample of 300 points was used as our quality data sample which would be used for training and verification.

2.2. Detect Cycles

For input of the model we decided to use power readings. We decided to go with power readings as one would expect a direct correlation between energy consumed and material removed. Also, it depends very little on parts of the machine which are not responsible for cutting (eg. housing) and thus we can make our system more flexible for various machines. Also, power readings can be easily obtained using conventional power meters.

As we plan to make an autonomous correlation model, we need to be able to correctly identify the start and end times of a production cycle. In relatively newer machines, this data can be obtained via the controller. However, the number of machines which are not capable of that are quite large. So we need to come up with a method to let the machine identify a cycle and record its start and end time. We use a simple pattern matching approach here. We are going to analyze the power readings we get as a stream of data and in it look for recognisable patterns which would indicate whether and where a cycle occurred. So to start, we need a known producing cycle which we can use as a template for comparisons. Once we decide this template we use an algorithm called **Dynamic Time Warping** (DTW).

2.2.1. *Dynamic Time Warping*

Dynamic time warping is a method to compute the similarity between two sequences. To compare these two sequences, one can *warp* (stretch/compress) the sequences non-linearly. Originally used as a tool to detect speech patterns in *Autonomous speech recognition*, DTW is now used in many applications which need to compare one sequence to another. One can limit the amount of warping and thus ensure that realistic results are obtained.

The reason why it is appropriate for us to use is that it gives us a very convenient method to compare two sequences and assign a distance between them. It's robust nature, which comes as a result of the warping allowed, lets us enter data with slight offsets along time, delayed peaks, noise, etc and obtain a good distance measure which can be effectively used to distinguish between a cycle and a non cycle.

2.3. *Extracting input from power consumption pattern*

Now we have the start and end times for 300 producing cycles and we can use this data to get power consumption patterns for each individual part and use it as necessary. However, we can't use these patterns directly as inputs for any machine learning algorithm as it has many inherent problems. The main reason being that small time offsets in the time series data could throw off the result by a large amount. So we basically needed to process our power pattern to extract more stable information from it. However in the process of making it stable, we also have to make sure that we don't lose the individuality of the pattern as we depend on this individual character for our quality reading. We resolved this using two methods:-

1. Scalar extractions from pattern.
2. Cycle splitting.

2.3.1. *Scalar extractions from splitting*

Here we extract from our cycle a set of scalars which are stable and do not get affected due to small changes like minor time scaling, peaks, etc. These scalars could be mean, standard deviation, DTW distance, median, mode, amplitude, etc. However, as pointed out earlier, using these scalars alone would render

many cycles to be same as we are wasting a lot of information. So we use, not one, but a combination of these scalars as now we get more options of variability thus giving us our desired result of stable and unique inputs.

2.3.2. *Cycle splitting*

This is a much more important process and increases the effectiveness of scalar extraction. What we essentially do cut the cycle into a number of parts using different methods (which will be discussed later) and apply the scalar extraction for each segment. So now we increase our number of inputs and also increase the individuality of each input in a stable way.

Much more importantly, when we later analyze our results, we can actually pinpoint which segments of the cycle are primarily responsible for the quality parameter being measured (in our case, R_a). This kind of information can prove to be really useful as now we can know what part of the cycle is primarily responsible for which quality parameter. This can in turn help us detecting tool damages when we observe a deviation in certain quality parameter. Splitting the cycle gives us one more important advantage. Now we can calculate the quality metric of a cycle by including one split at a time, into our input set, sequentially. This can be really helpful in predicting parameter qualities before the machining is over and thus can help us in taking corrective action.

We primarily employed four different methods to split cycles. They were:-

1. **Uniform splitting:** As the name suggests, we split the cycle into equally sized sections. The number of parts would depend on how much accuracy we want and how much computation cost we can afford.
2. **Forward feedback splitting:** This is a process in which we make a cut by splitting a section into two equal halves and decide whether we keep that cut by checking if the amount of improvement we get is significant enough. It is slightly biased to create more cuts in the beginning of the cycle as we move on to the former section after a cut has been made and move to the next section only if no more cuts can be made in the former one.
3. **Reverse feedback splitting:** This is almost like the Forward Feedback Splitting but instead of biasing cuts to the start, this one biases cuts to the end.
4. **Phase dependent cuts:** These are cuts which signify a change in process in the machining. Theoretically, these should be the most optimum cuts as they point towards specific process. However, finding them out autonomously can be tricky if wanted for all situations thus these cuts have to be entered manually.

Now, after splitting the cycle appropriately and using scalar extraction for each split, we have a bunch of variables for a cycle which we will later use to get a wider set of data.

2.4. Model Selection

There are a plethora of machine learning algorithms which we could use to establish a relation between our input and output. We decided to try the very simple models and thus decided to use *Polynomial Regression*. However, as our number of variables were bound to be high, we decided to keep the model complexity to a minimum and thus decided to go forward with *Regularised Polynomial Regression*.

2.4.1. Regularised Regression

Say we try to create a linear relation between input variables $x_1, x_2, x_3, \dots, x_n$ and output variable y . As it is linear, we can assume our equation to be

$$y = a_0 + a_1x_1 + a_2x_2 + a_3x_3 + \dots + a_nx_n \quad (1)$$

Where $a_1, a_2, a_3, \dots, a_n$ are the coefficients which we have to find.

Let us define input variable \mathbf{X} and \mathbf{A} as

$$\mathbf{X}_i = (x_1^i, x_2^i, x_3^i, \dots, x_n^i) \quad (2)$$

$$\mathbf{A} = (a_1, a_2, a_3, \dots, a_n) \quad (3)$$

where x_j^i is the i^{th} instance of the j^{th} variable. Thus, we get

$$y_i = a_0 + a_1x_1^i + a_2x_2^i + a_3x_3^i + \dots + a_nx_n^i \quad (4)$$

Now, say we have the p values of \mathbf{X}_i and its corresponding output which we denote by \mathbf{Y}_i . So to find the best approximation of \mathbf{A} , we minimize the difference between \mathbf{Y}_i and y_i . Let our final coefficients' set be \mathbf{A}_{LS} . Using least square method, we get

$$\mathbf{A}_{LS} = \arg \min_{\mathbf{A}} \left(\sum_{i=1}^p (y_i - \mathbf{Y}_i)^2 \right) \quad (5)$$

which can be re-written as

$$\mathbf{A}_{LS} = \arg \min_{\mathbf{A}} \left(\sum_{i=1}^p (\mathbf{A}^T \mathbf{X}_i - \mathbf{Y}_i)^2 \right) \quad (6)$$

Thus as a result of linear regression, we finally get

$$y = \mathbf{A}_{LS}^T \mathbf{X} \quad (7)$$

when x_j are polynomial forms of each other, this is called polynomial regression.

However, when the number of variables are high, i.e., when n is large, we tend to get highly complex models with high values of a_i which tend to overfit

y as a function of \mathbf{X} . To remedy this, we include a regularization parameter λ , such that, $\lambda > 0$ and

$$\mathbf{A}_{LS} = \arg \min \left(\sum_{i=1}^p (\mathbf{A}^T \mathbf{X}_i - \mathbf{Y}_i)^2 + \lambda \sum_{j=1}^n a_j^2 \right) \quad (8)$$

Thus, now we can control, via λ , how complex a system we want. Too high a λ gives us an overfit solution while too low a λ gives us an underfit one.

2.5. Model Creation

To implement *Regularised Polynomial Regression*, we first need to use our present variable set (which we got after extracting scalars from split sections of a cycle) to form polynomial variables. Firstly we divide our 300 cycles into two sets, a *training set* and a *testing set*. each cycle is assigned to one of these sets randomly and the ratio between the number of elements in the two sets can be decided by us. we chose a 1:1 ratio. Now, in the training set, for each variable (scalar from a split section), we have 150 other instances of the same corresponding to 150 cycles. This gives us a vector of that corresponding variable. This way we have a vector corresponding to all the variables. These vectors will be used to create the polynomial variables which will eventually be used as inputs into the *Regularised Linear Regression*. This we do in two ways:-

1. **Orthogonal polynomials-** We create new vectors which are orthonormal to each of the vector of variables we have. The degree of these orthogonal variables were varied to see which would yield reasonable results.
2. **Mixed polynomials-** Here we create new vectors using the old variable vectors (denoted as $x_1, x_2, x_3, \dots, x_n$) which can be shown asllll

$$x_{new}(i_1, i_2, i_3, \dots, i_n) = \prod_{j=1}^n x_j^{i_j}; \quad (9)$$

where

$$i_1, i_2, i_3, \dots, i_n = 0, 1, 2, 3, 4, \dots; \quad (10)$$

and

$$\sum_{j=1}^n i_j \leq p; \quad (11)$$

Where p is the maximum degree which we vary to get optimum results. As expected, we get more number of variables using this method instead of the orthogonal polynomials one. However, it helps us know if there are any meaningful interactions within the variables.

Now that we have our input vectors' set created, we just a *Regularised Linear Regression* to get a list of possible relations between our input quality. In the next section we discuss how we select the right parameters from this given list.

2.6. Parameter comparison and selection

The output of the *regularised polynomial regression* is a list of solutions. Each solution contains a collection of coefficients which we use for a weighted summation of our newly calculated variable vectors. To compare between these various solutions, we use them and predict quality parameter values for our test set and then we compare these predicted values against our polynomial fit values (from 2.1) and find out the *root mean square error (RMSE)*, where we assume that the polynomial fit values are the true values. We use this *RMSE* as a metric to compare between the solutions and pick the solution with minimum *RMSE*.

AppendixA. Appendix name

Appendix, only when needed.

You can use either BibT_EX:

Or plain bibliography:

- [1] Frank Mittelbach and Michel Goossens: *The L^AT_EX Companion Second Edition*. Addison-Wesley, 2004.
- [2] Scott Pakin. The comprehensive L^AT_EX symbol list, 2005.