# USING FUNCTIONAL DATA ANALYSIS FOR TIME-DEPENDENT OPTIMIZATION OF BATCH PROCESSES

Hadley Myers, JMP

**NCS** Non Clinical Statistics Conference

**October 3-4-5 2018** Maison de la Chimie, Paris

- Wolfgang Ketterle, Nobel Lecture (2001)
  - "Imagine how many aspects of nature we would miss if we lived on the surface of the sun…without refrigerators."



- If we can create conditions that haven't been created before, then we'll make new discoveries.
- If we look at things in a way that hasn't been done before, then we'll see new things.
- If we analyze data using new methods, then we'll gain new insights.

# FUNCTIONAL DATA ANALYSIS

- Very often, data will be "telemetric" in nature - many repeated measures of several metrics through time.
  - This is true of data from many sources.
    - Machines output
    - Traditional time series data
    - Sensor data
    - Vibration signals
- A wide variety of specific tools and methods have been created to deal with this type of data:
  - Signal Processing
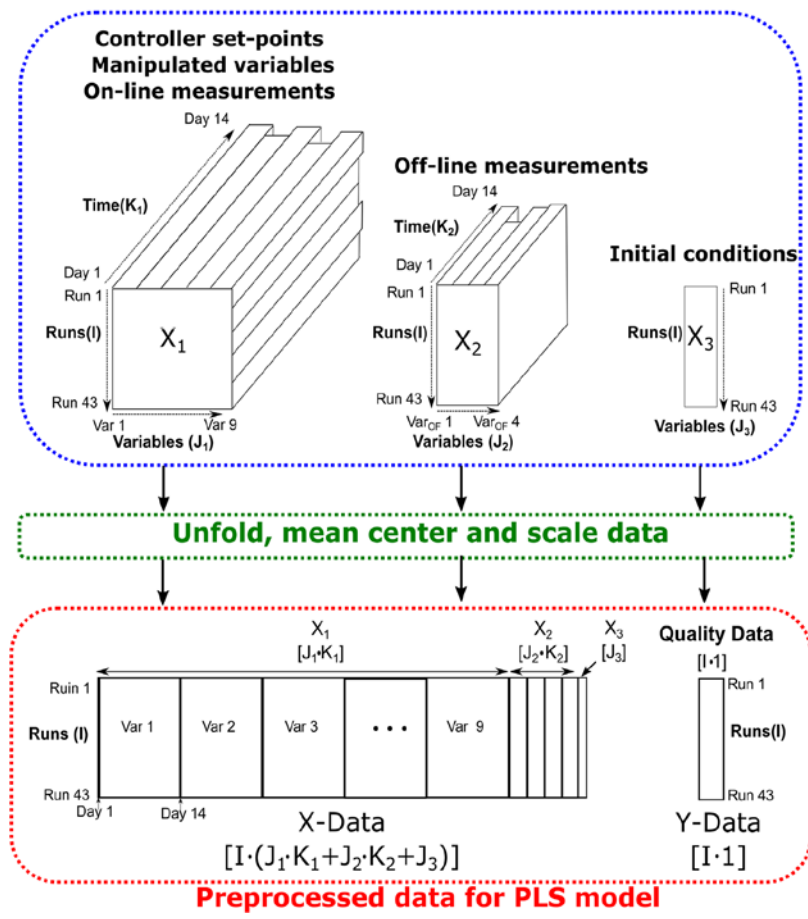  - ARIMA Time Series
  - Partial Least Squares
  - Growth Curves via SEM
  - Mixed Models

# FUNCTIONAL DATA ANALYSIS

- Many products are made in batches by machines that now have many sensors embedded in them
- Sensors record things like temperature, pressure, feed rate, chemical content (ammonia, $CO_2$, ethanol, sugar), vibration, etc.
- Companies care about end results:
  - Yield: the quantity of product created (yield)
  - Quality: Measurable properties of the product (flavor, room temperature viscosity, shear strength, chemical composition)
- They want to understand how the sensor readings relate to the end results
  - To fix 'bad batches', or terminate their production early
  - Reduce occurrence of bad batches (process improvement)
- This is not a new problem - Due to the explosion of data access a lot more people want to take advantage of functional data

# FUNCTIONAL DATA ANALYSIS

- Traditional approaches are too often inadequate and overcomplicated
  - Converting data to wide format (one input variable per time period) and using PLS.
    - Data cleaning step can be very time-consuming.
    - Sparse table if time-points not aligned, lost data if sizes not equal.
    - Difficult to interpret results for optimization (may be possible for early flagging of batches)
  - Least Squares modeling of summary statistics (mean, min, max, etc.).
    - Too simple, all time-dependent information is lost.
    - Model says nothing about the shape of the functions.
  - Fitting logistic curves, using parameter estimates as features.
    - Very limited in the set of shapes of curves that can be fit.
    - More flexibility needed than simple logistic curves and Gaussian peak models (again, too simple).

# DEALING WITH TELEMETRIC DATA
## (TRADITIONAL APPROACH)



Goldrick et al.: MVDA of Trisulfide Bond Formation, Biotechnology and Bioengineering, Vol. 114, No. 10, October, 2017

# USES OF FUNCTIONAL DATA ANALYSIS

1. Functional factors, constant responses.
2. Constant factors, functional responses.
3. Both functional factors and functional responses.

Example: Maximizing the yield of human insulin produced by modified yeast cells.

# TABLE OF RAW DATA



- There are 100 batches in total
- 100 time measurements per batch
- Measurements were taken at fixed time intervals
- This isn't always the case!

Plotting the factors over time for each batch reveals the complexity of the problem.

Where(BatchID = 3)

Exploring batches individually further emphasizes the challenge ahead.

# GOOD BATCHES VS BAD BATCHES

What we'll find:
1) We can achieve a 74% Yield.
2) Ethanol, Molasses, NH3 and Air are significant.
3) These are their ideal profiles:



We could try to look for characteristics of good vs bad batches, but what defines "good?"
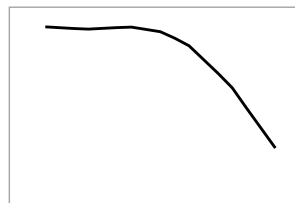
Can "good" be "better?"
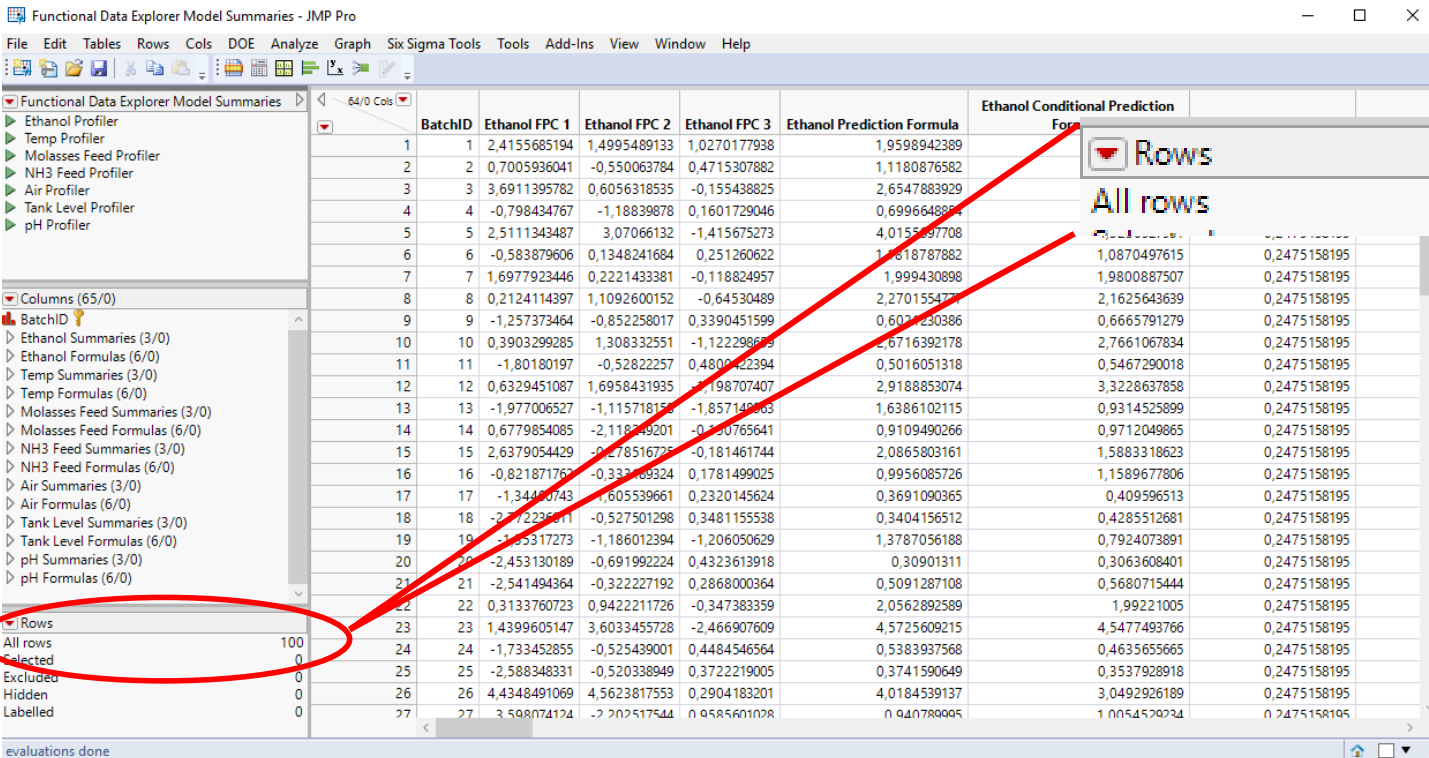
# FUNCTIONAL DATA ANALYSIS



Functional Data Analysis modeling for determining Functional Principal Components to be used in prediction.

# SAVING FPC AND EIGENFUNCTIONS FOR EACH BATCH



- We now have 100 functional summaries (one for each batch)

# MODELING WITH FPCS AS INPUTS

**Actual by Predicted plot**



| | |
|---|---|
| RSquare | 0,733979 |
| RSquare Adj | 0,731265 |
| Root Mean Square Error | 0,020353 |
| Mean of Response | 0,53436 |
| Observations (or Sum Wgts) | 100 |

**Generalized Regression for Final Yield**
**Adaptive Lasso with AICc Validation**
**Solution Path**



**Parameter Estimates for Original Predictors**

| Term | Estimate | Std Error | Wald ChiSquare | Prob > ChiSquare | Lower 95% | Upper 95% |
|---|---|---|---|---|---|---|
| Intercept | 0,53436 | 0,0020203 | 69956,999 | <,0001* | 0,5304003 | 0,5383197 |
| NH3 Feed FPC 2 | -0,000173 | 2,4377e-5 | 50,628148 | <,0001* | -0,000221 | -0,000126 |
| Ethanol FPC 2 | -0,012168 | 0,0018192 | 44,736472 | <,0001* | -0,015733 | -0,008602 |
| Molasses Feed FPC 2 | 0,0000244 | 4,0897e-6 | 35,568946 | <,0001* | 1,6375e-5 | 0,0000324 |
| Air FPC 1 | -2,63e-6 | 1,1457e-6 | 5,2684672 | 0,0217* | -4,875e-6 | -3,842e-7 |
| Ethanol FPC 1 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Ethanol FPC 3 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Temp FPC 1 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Temp FPC 2 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Temp FPC 3 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Molasses Feed FPC 1 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| Molasses Feed FPC 3 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| NH3 Feed FPC 1 | 0 | 0 | 0 | 1,0000 | 0 | 0 |
| NH3 Feed FPC 3 | 0 | 0 | 0 | 1,0000 | 0 | 0 |

- Significant factors

| Term | Estimate | Std Error | Wald ChiSquare | Prob > ChiSquare |
|---|---|---|---|---|
| Intercept | 0,53436 | 0,0020203 | 69956,999 | <,0001* |
| NH3 Feed FPC 2 | -0,000173 | 2,4377e-5 | 50,628148 | <,0001* |
| Ethanol FPC 2 | -0,012168 | 0,0018192 | 44,736472 | <,0001* |
| Molasses Feed FPC 2 | 0,0000244 | 4,0897e-6 | 35,568946 | <,0001* |
| Air FPC 1 | -2,63e-6 | 1,1457e-6 | 5,2684672 | 0,0217* |

# INTERPRETATION OF RESULTS



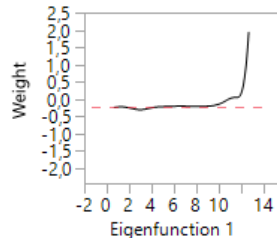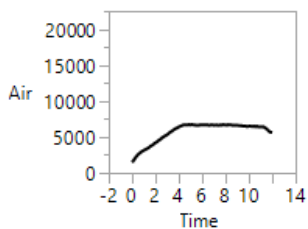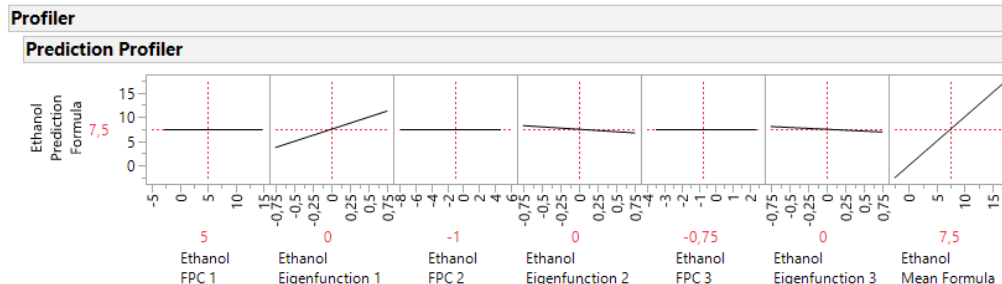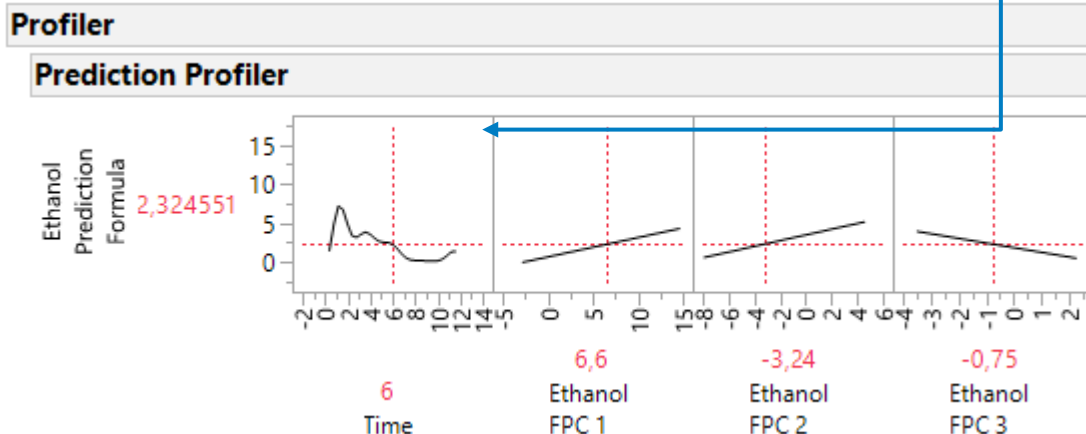| Term | Estimate | Std Error | Wald ChiSquare | Prob > ChiSquare |
|---|---|---|---|---|
| Intercept | 0,53436 | 0,0020203 | 69956,999 | <,0001* |
| NH3 Feed FPC 2 | -0,000173 | 2,4377e-5 | 50,628148 | <,0001* |
| Ethanol FPC 2 | -0,012168 | 0,0018192 | 44,736472 | <,0001* |
| Molasses Feed FPC 2 | 0,0000244 | 4,0897e-6 | 35,568946 | <,0001* |
| Air FPC 1 | -2,63e-6 | 1,1457e-6 | 5,2684672 | 0,0217* |

# INTERPRETATION OF ONE FACTOR (ETHANOL)

Predictive formula for Ethanol, from FDA B-spline fitting



- Ethanol FPC 1 • Ethanol Eigenfunction 1

+

- Ethanol FPC 2 • Ethanol Eigenfunction 2

+

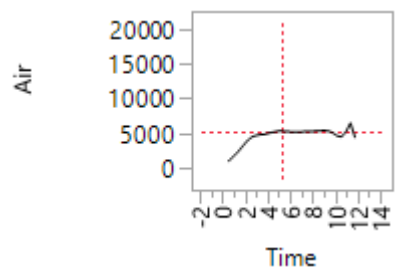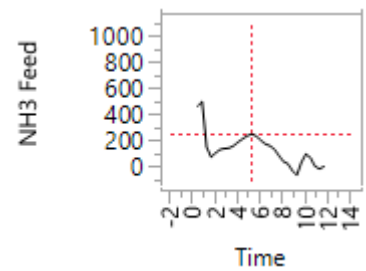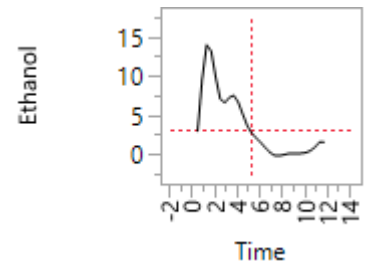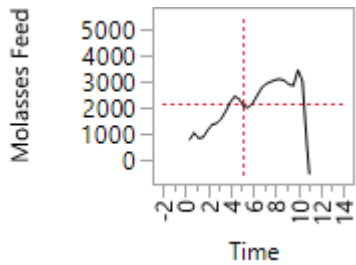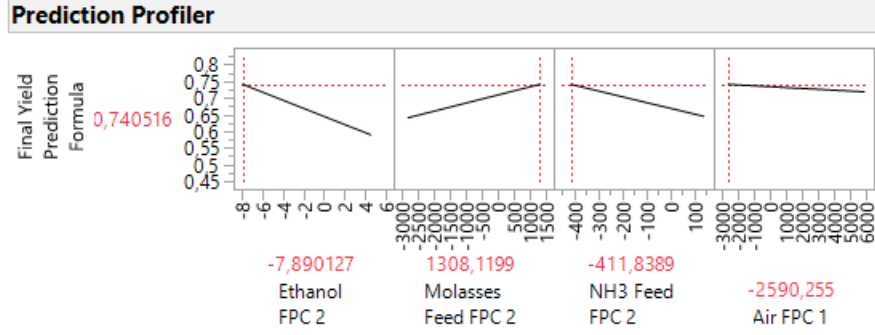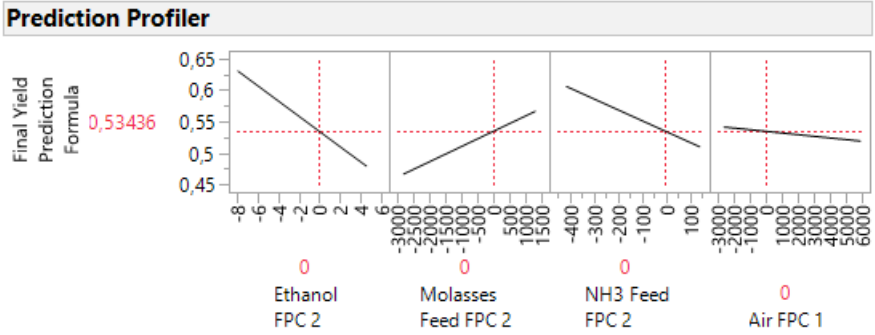- Ethanol FPC 3 • Ethanol Eigenfunction 3

+

Ethanol Mean Formula

These are functions of time, which can be consolidated.

The Prediction Profiler is used to visualize and explore predictive models.

# OPTIMAL YIELD RECIPE

Optimizing the values of the FPCs to maximize Yield, we have everything needed to find the optimal profiles.



These profiles are predicted to result in a Yield of 74%.

- Using Functional Data Analysis over traditional methods, we can:
  - Dramatically reduce total time to meaningful results,
  - Utilize all data while preserving time-dependent info,
  - Generate more interpretable knowledge output,
  - Better engage with subject-matter experts.

- Question: Can you think of opportunities within your organization where you could apply this method?