



# DataHow Coding Challenge

Presented by Adriana Mohap

31 October 2024



# Task Description

1. Perform any data analysis you consider relevant.
2. Build any data processing pipeline and model of your choice.
3. Evaluate the performance of your model.
4. Prepare your code so that it can generate predictions and results on a holdout test set (which we will provide before the interview).



# Data Description

- X:{name}: Represents measurements of process conditions, which are variables inherent to the process.
- W:{name}: Represents control conditions, which are measurements of the controlled parameters in the process.
- Z:{name}: Represents control setpoints, which are operator-defined control conditions that remain constant throughout the process.
- Y:{name}: Represents process attributes, which are measured at the end of the process.



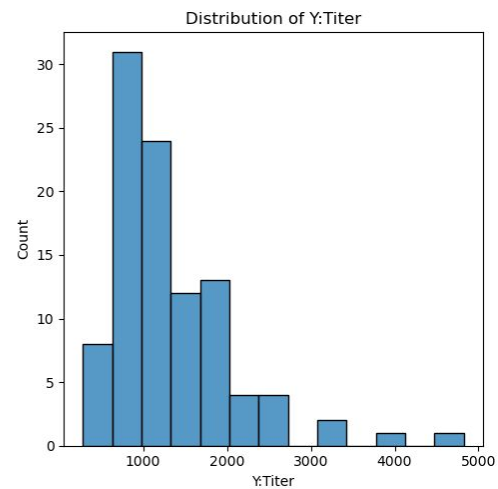
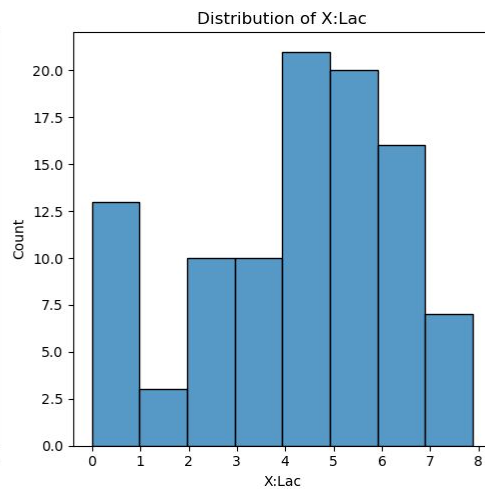
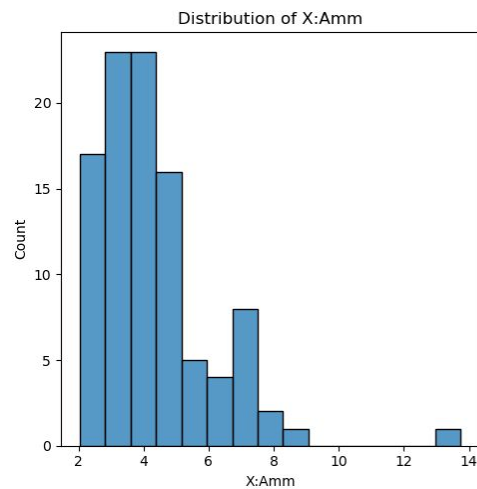
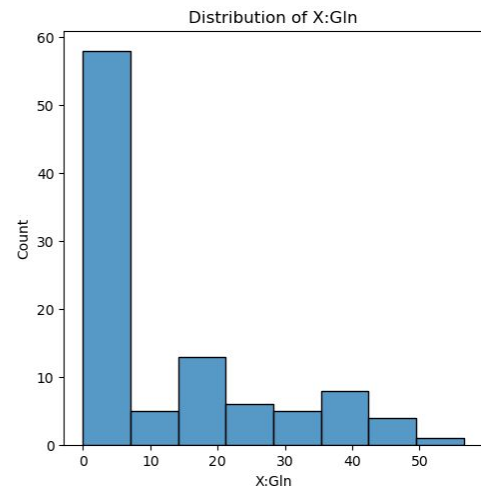
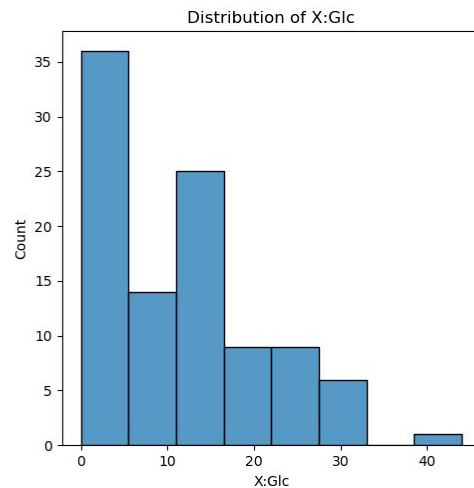
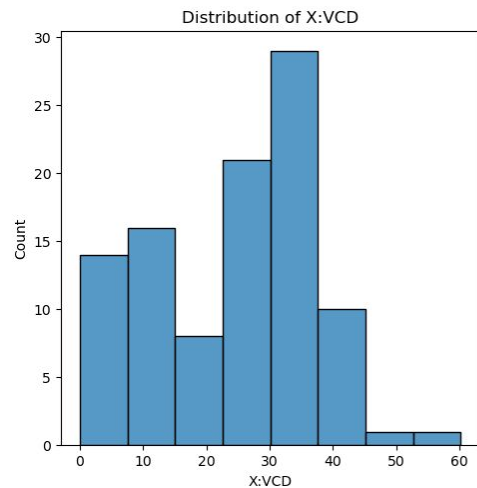
## Task 1: Data Analysis

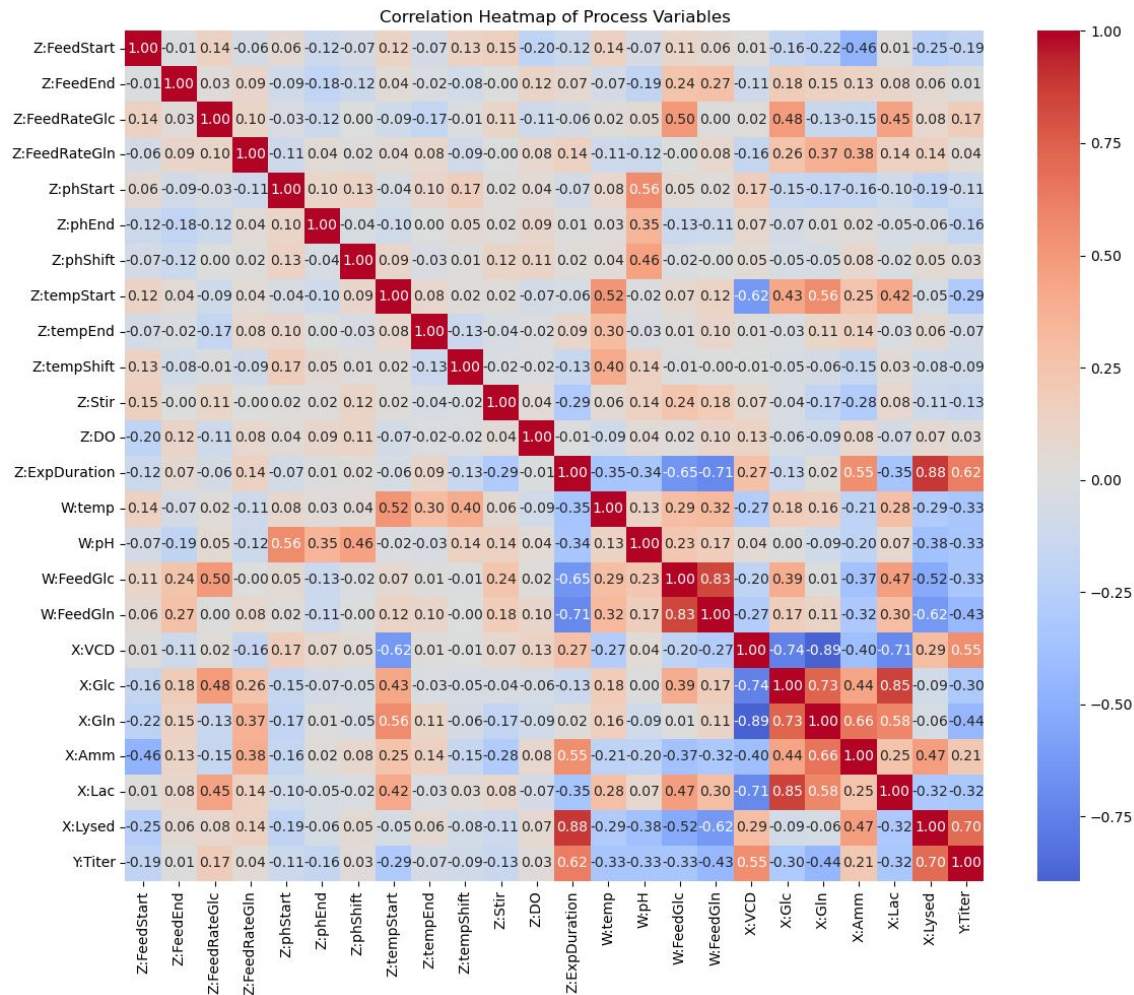
- Load and preprocess data
- Aggregate descriptive statistics of all variables: mean, standard deviation, min and max values
- Plot the process variable distributions against the target value
- Analyze correlations between all variables

Process Variables (X:):						
	X:VCD	X:Glc	X:Gln	X:Amm	X:Lac	X:Lysed
count	990.000000	990.000000	990.000000	990.000000	990.000000	990.000000
mean	13.425586	7.210130	7.406499	2.023753	3.271055	0.028822
std	12.377215	6.399089	9.829110	1.809282	2.000817	0.061393
min	0.003421	0.000000	0.000000	0.000000	0.000000	0.000000
25%	3.704009	2.910059	0.147580	0.515444	1.580819	0.000000
50%	8.991166	5.133313	4.136096	1.602359	3.488986	0.007234
75%	21.486081	9.738994	9.719193	3.103277	4.843664	0.024610
max	61.242464	44.011310	62.089637	13.744255	8.078131	0.525695

Control Variables (W:):				
	W:temp	W:pH	W:FeedGlc	W:FeedGln
count	990.000000	990.000000	990.000000	990.000000
mean	36.883461	6.941144	2.963942	5.250168
std	0.649905	0.326060	1.987881	3.084908
min	35.090909	6.025253	0.000000	0.000000
25%	36.414141	6.696970	0.000000	0.000000
50%	36.878788	6.929293	3.272727	6.696970
75%	37.404040	7.212121	4.606061	7.343434
max	37.989899	7.494949	5.979798	7.989899

Setpoint Variables (Z:):													
	Z:FeedStart	Z:FeedEnd	Z:FeedRateGlc	Z:FeedRateGln	Z:phStart	Z:phEnd	Z:phShift	Z:tempStart	Z:tempEnd	Z:tempShift	Z:Stir	Z:D0	Z:ExpDuration
count	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000	100.000000
mean	2.000000	11.000000	4.000000	7.000000	7.000000	6.500000	10.000000	37.000000	36.000000	10.000000	200.000000	55.000000	8.900000
std	0.710669	1.214392	1.154642	0.577321	0.288660	0.288660	2.348436	0.577321	0.577321	2.348436	28.866041	14.433020	2.032563
min	1.000000	9.000000	2.020202	6.010101	6.505051	6.005051	6.000000	36.010101	35.010101	6.000000	150.505051	30.252525	7.000000
25%	1.750000	10.000000	3.020202	6.510101	6.755051	6.255051	8.000000	36.510101	35.510101	8.000000	175.505051	42.752525	7.000000
50%	2.000000	11.000000	4.000000	7.000000	7.000000	6.500000	10.000000	37.000000	36.000000	10.000000	200.000000	55.000000	8.500000
75%	2.250000	12.000000	4.979798	7.489899	7.244949	6.744949	12.000000	37.489899	36.489899	12.000000	224.494949	67.247475	10.000000
max	3.000000	13.000000	5.979798	7.989899	7.494949	6.994949	14.000000	37.989899	36.989899	14.000000	249.494949	79.747475	14.000000







## Task 2: Build Data Processing Pipeline and Model

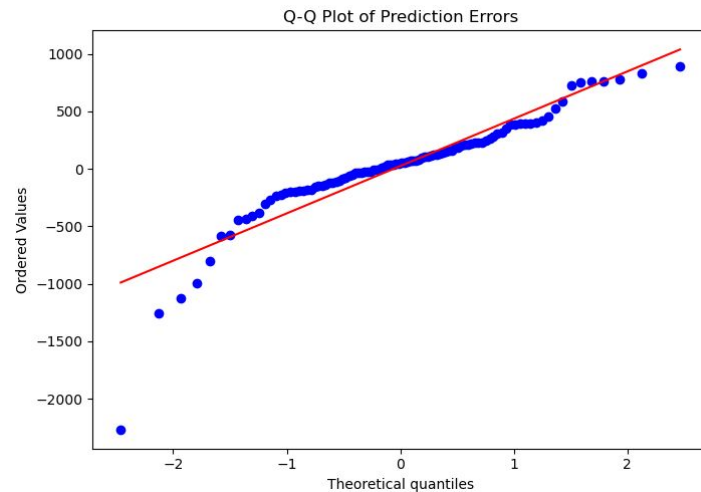
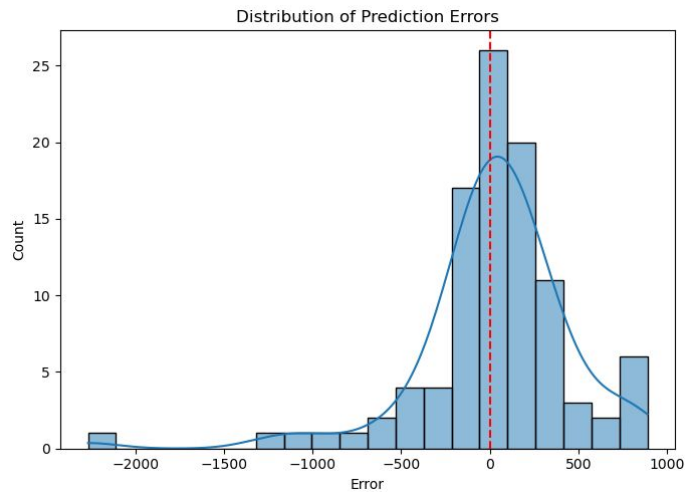
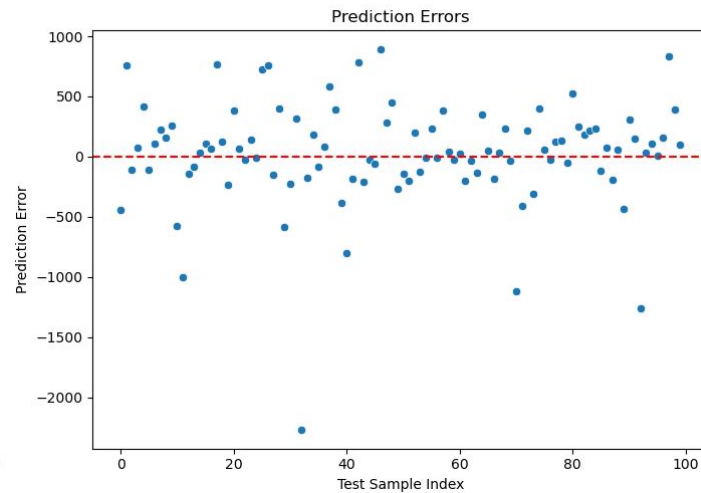
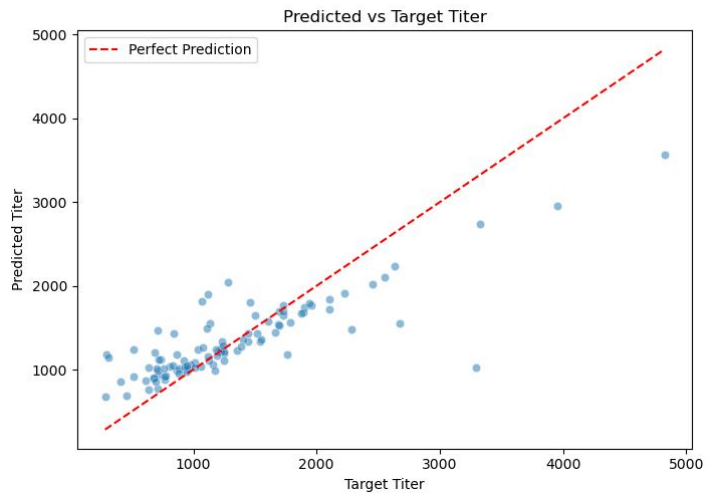
- Build two functions for feature engineering:
  - `engineer_features(process_data)`: **X variables**: mean, std, min, max; **W variables**: mean only; **Z variables**: initial values only
  - `engineer_features_baseline(process_data)`: Baseline feature engineering using **only mean values of X** variables
- Random forest regression as model
  - no assumptions about the underlying distribution of the data
  - works well with high dimensional data and collinearity in features



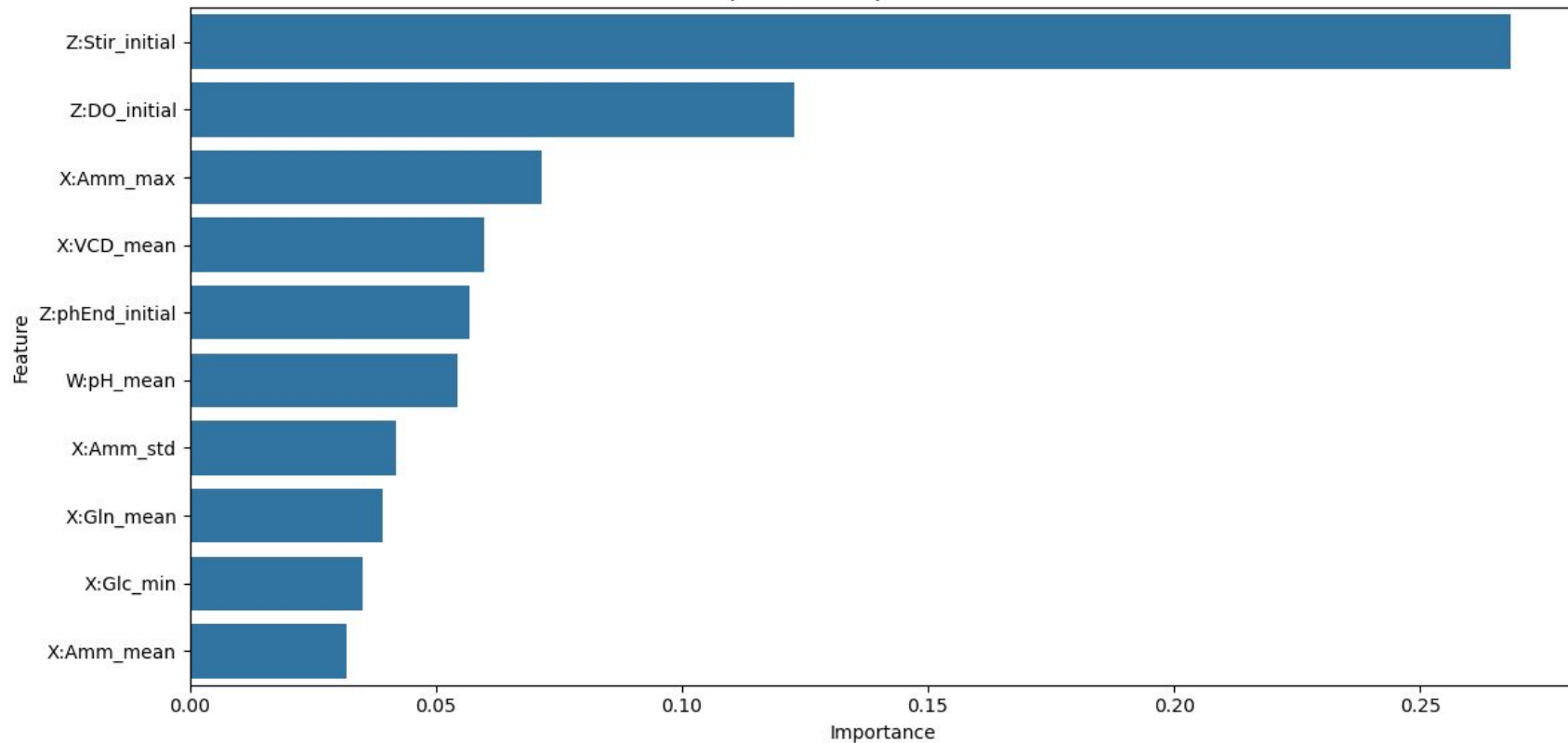


## Task 3: Evaluate Performance

- Metrics over the Training and Test set:
  - R Squared
  - RMSE
  - MAE
  - MAPE
- Visual Inspection of the predictions and the residuals



Top 10 Most Important Features





## Future Work

- Gain domain knowledge about the process variables in order to adjust the selected features
- Create features that accurately reflect their relationships (interaction terms, growth rates etc.)
- Drop features that show high correlations
- Perform hyperparameter tuning; Cross-validation
- Use time variable as separate variable
- Select a model, which is able to capture time series data more accurately (for example LSTM or Recurrent networks) or use ensemble techniques