# Data Drive Model for Outlier Detection

**Melanie Santiago**

Industry Pricing Branch Chief

Producer Price Index

U.S. Bureau of Labor Statistics

39th Voorburg Group Meeting
September 22 – 26, 2025

# Outlier detection in the US PPI

## Use of tolerance thresholds

- Quality checks on reported microdata
- System flag for review
- Price verification by BLS economist
- Justification for price change logged

## Traditional method for setting tolerance level

- Set when new industry sample introduced
- Values for positive and negative thresholds
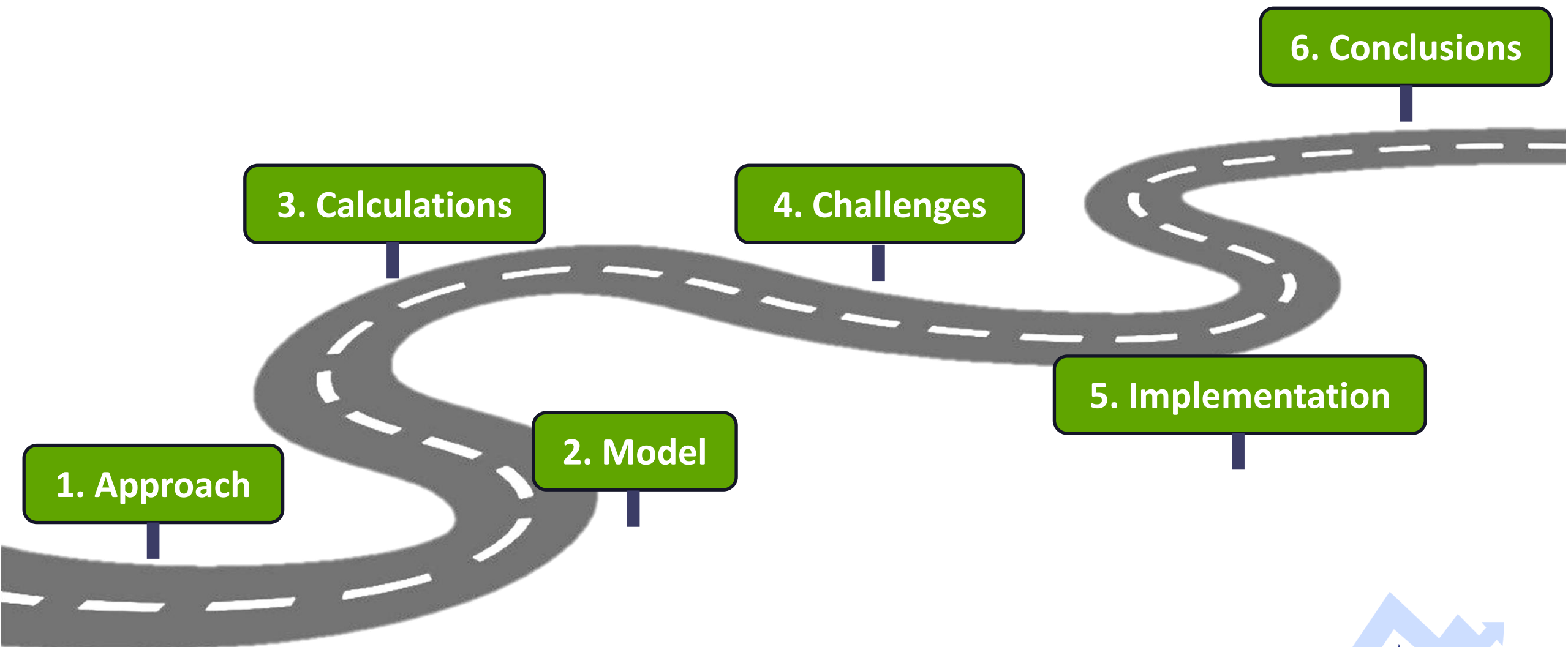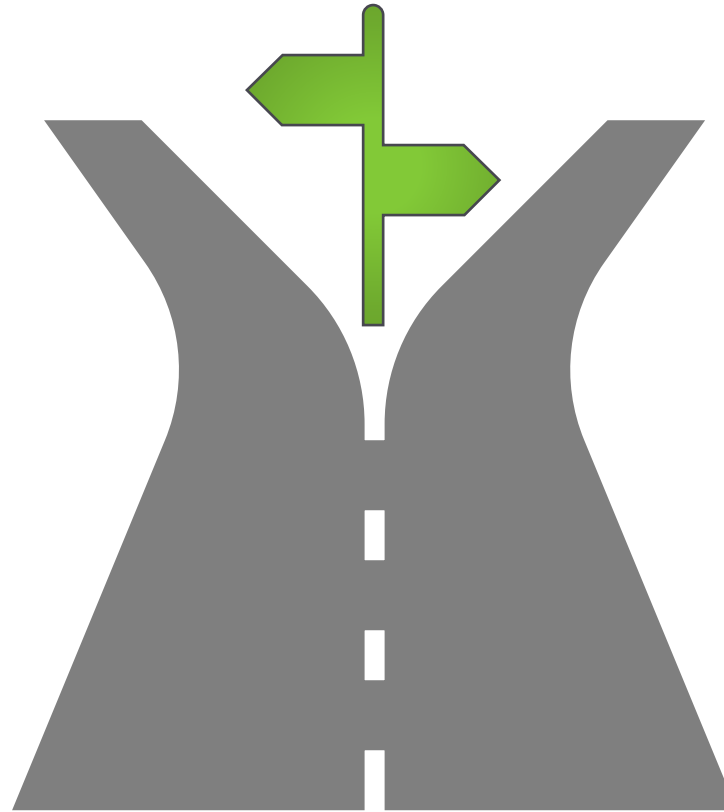- Static values in data processing system

# Goals



- Data-driven decision making
- Improved outlier detection
- Dynamic process
- Operational efficiencies

# Roadmap



6. Conclusions

3. Calculations

4. Challenges

5. Implementation

2. Model

1. Approach

# Approach

Develop ARIMA Model: Autoregressive Integrated Moving Average

Explore alternative modeling techniques

BLS

# Modeling techniques

| | Accuracy | Generalizability | Computational Efficiency | Interpretability |
|---|---|---|---|---|
| **ARIMA: Autoregressive Integrated Moving Average** | ✓ | ✓ | ✓ | ✓ |
| **ETS: Error, Trend and Seasonality** | ✗ | | | |
| **Machine Learning** | | | | ✗ |
| **Deep Learning** | | | ✗ | |

# Model

$$ARIMA\ (p, d, q)(P, D, Q)[m]$$

Parameters

p: range of values allowed for autoregressive terms
d: range of values allowed for order difference
q: range of values allowed for moving average terms

P: range of values allowed for seasonal autoregressive terms
D: range of values allowed for seasonal order of differences
Q: range of values allowed for seasonal moving avg terms

m: seasonality

# Seasonal model

$$ARIMA(start\_p = 0,\ max\_p = 23,\quad d = None,\ max\_d = 1,$$
$$start\_q = 1,\ max\_q = 23,\quad max\_P = 2,\ D = 0,\ max\_Q = 1,$$
$$m = 12,\ seasonal = TRUE,\quad stepwise = TRUE$$

- Up to 24 months of net price changes
- Up to 24 months of past forecast errors used to predict future errors
- Incorporate seasonal trends from up to 2 prior 12-month periods

BLS

# Non-seasonal model

$$ARIMA\ (start\_p = 0,\ max\_p = 11,\quad d = None,\ max\_d = 1,$$
$$start\_q = 1,\ max\_q = 11,\quad seasonal = FALSE,$$
$$stepwise = \quad TRUE$$

- Up to 12 months of net price changes
- Up to 12 months of past forecast errors used to predict future errors
- No seasonal trends used

# Data preparation

Consolidate price change data into two price change values per month

| Median of top 10 percent of all positive net price changes in each industry – for positive outlier detection model | Median of bottom 10 percent of all negative net price changes in each industry – for negative outlier detection model |

# Calculation

1. Enter median of upper 10% of positive or lower 10% of negative net price changes into Auto-ARIMA

2. Select seasonal or non-seasonal model based on amount of historical data available

3. Use Auto-ARIMA to generate parameter options

4. Determine optimal parameters by AIC for ARIMA model for tolerance level predictions

5. Create error dataset (difference between predicted and actual value for all historical time periods)

6. Determine absolute value of median error

7. Create deviation dataset (median of differences between Median Error and actual error for each time period)

8. Calculate tolerance level for a given time period

# Challenges

## Structural changes (Combined or recoded industries)

Model requires two years of price data for each industry code

New industry codes have limited historical data

Manual process to link data from the prior and component indexes to new and combined industry codes
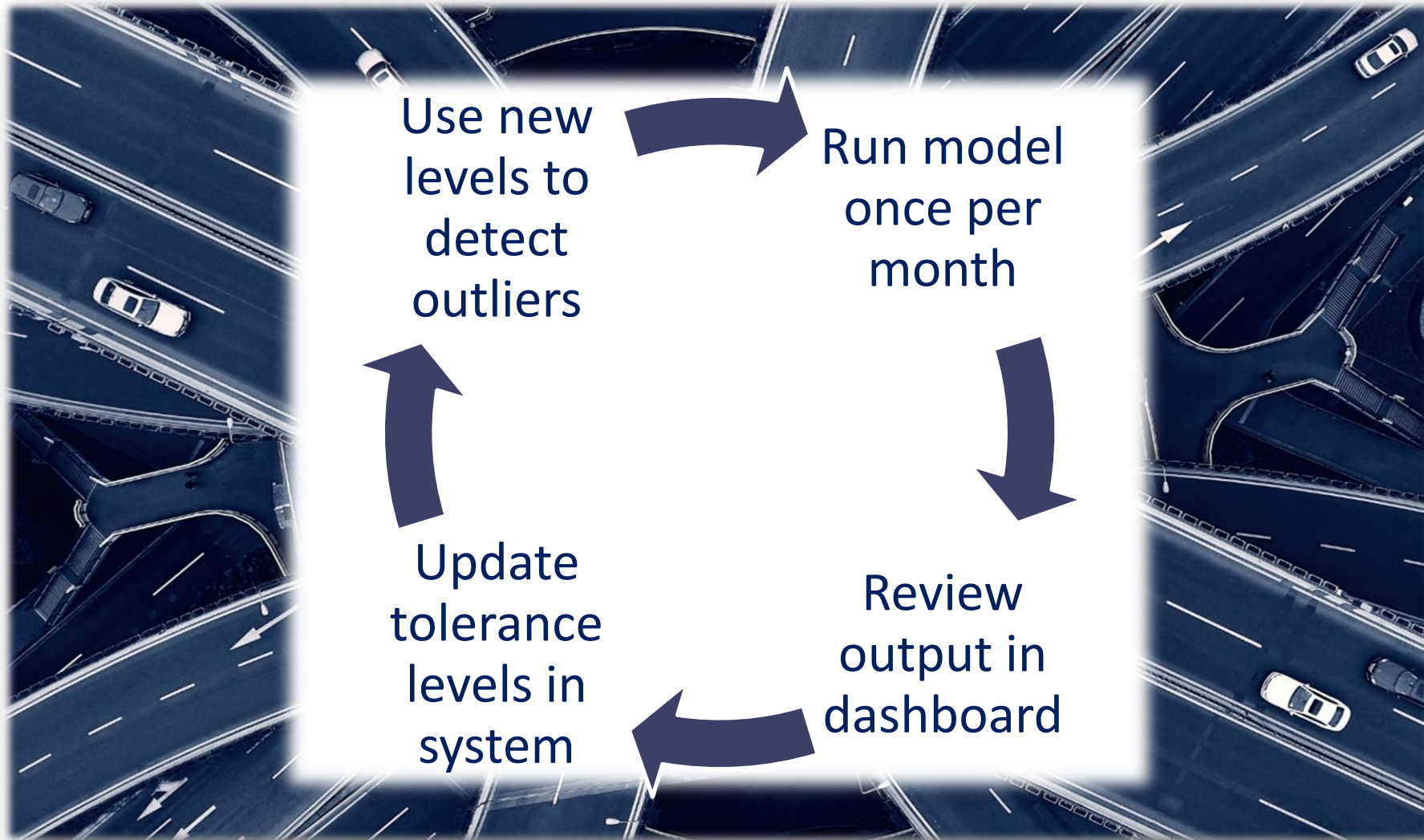
## Product level tolerances

Different products produced within an industry may exhibit different price trends

Lower-level indexes are prone to more frequent structural changes

Structural changes result in a lack of required historical data

# Implementation



Use new levels to detect outliers → Run model once per month → Review output in dashboard → Update tolerance levels in system →

BLS

# Dashboard

# Conclusions

■ Data driven decision making

▶ Microdata review efficiency

▶ Reduce workloads?

▶ Improve resource allocation



You have arrived at your destination

# Contact Information

**Melanie Santiago**
Industry Pricing Branch Chief
Producer Price Index
U.S. Bureau of Labor Statistics
www.bls.gov/ppi

202-691-7844

santiago.melanie@bls.gov

BLS