# Cyclic Boosting

A Pure-Python, Explainable, and Efficient ML Method

Felix Wick, June 2024

# This is old stuff.

# But it just works.

# Cyclic Boosting ML Methods

family of off-the-shelf, general-purpose supervised machine learning methods for both regression and classification tasks (focus on structured data)

closest relatives: Generalized Additive Models (not a deep learning approach)

main difference: estimation of factors for each bin of the different features (instead of estimation of parameters like coefficients in linear regression or weights in neural networks)

→ individual explainability

scientific papers describing the methods: **Cyclic Boosting**, **Demand Forecasting with Cyclic Boosting**

# Cyclic Boosting Library

**scikit-learn-like usage of library**

**open source: https://github.com/Blue-Yonder-OSS/cyclic-boosting**

**Python package: pypi**

**documentation: readthedocs**

# Cyclic Boosting | Different Modes/Scenarios

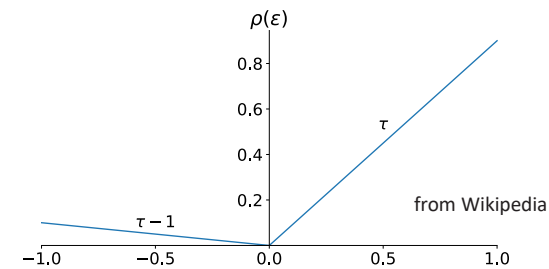| Mulitplicative Regression | Additive Regression | Classification | Negative Binomial Width | Exponential Price Elasticity | Background Subtraction |
|---|---|---|---|---|---|
| (conditional mean) | (conditional mean) | (probability) | (dispersion parameter) | (elasticity parameter) | (conditional mean) |
| $Y \in [0, \infty)$ | $Y \in (-\infty, \infty)$ | $Y \in [0, 1]$ | $Y \in [0, 1]$ | $Y \in [0, \infty)$ | $Y \in (-\infty, \infty)$ |
| Poisson / Negative Binomial distribution (link function *ln*) | Gaussian distribution (link function *identity*) | Bernoulli distribution (link function *logit*) | Negative Binomial distribution (link function *logit*) | exponential distribution (link function *ln*) | Gaussian distribution (link function *identity*) |
| *example* | *example* | *example* | *example* | *example* | *example* |
| demand forecasts (mean) | profit predictions | churn probability | demand forecasts as full probability distributions | individual price-demand elasticities | individual causal effects, e.g., customer targeting |

# Also Possible: Quantile Regression

quantile regression: estimate quantile $\tau$ of distribution instead of conditional mean by minimizing pinball loss

$$(1 - \tau) \sum_{y_i < \hat{q}_i} (\hat{q}_i - y_i) + \tau \sum_{y_i \geq \hat{q}_i} (y_i - \hat{q}_i)$$

instead of squared error loss (choice of loss function defines point estimate)

from Wikipedia

possible with various ML methods, including neural networks, tree-based methods (like random forests or gradient boosting), and Cyclic Boosting (minimized in each feature bin)
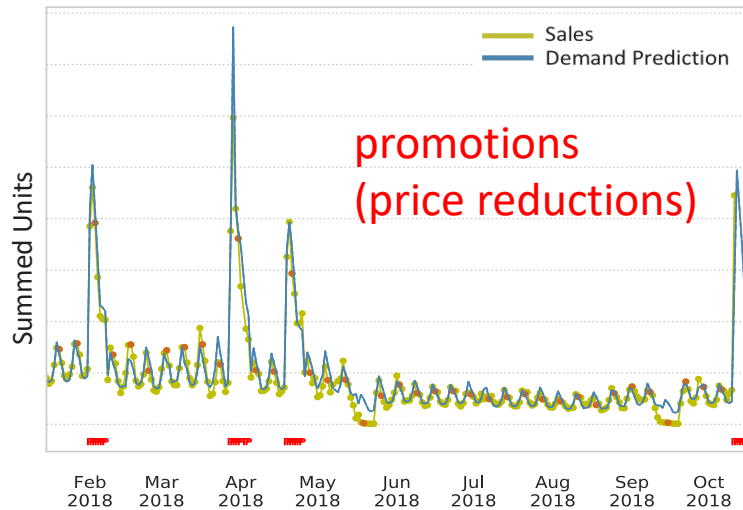
also: generic loss mode (e.g., maximum likelihood)

# Example Use Case: Demand Forecasting

**many individual time series to consider**

typical retail grocery chain:
- products (items): ~20k
- locations (stores): ~500
- daily/hourly aggregated sales



**advantages of machine learning over traditional univariate time series forecasting**

**combined learning on all time series** of product-location combinations (rather than separately optimizing individual time series)
→ **reduces variance** by exploiting commonalities

**natural consideration of many exogenous variables** (prices, promotions, holidays, weather, …)
→ **reduces bias**

to be noted:
- categorical features important (products and locations → high cardinality)
- mainly multiplicative effects
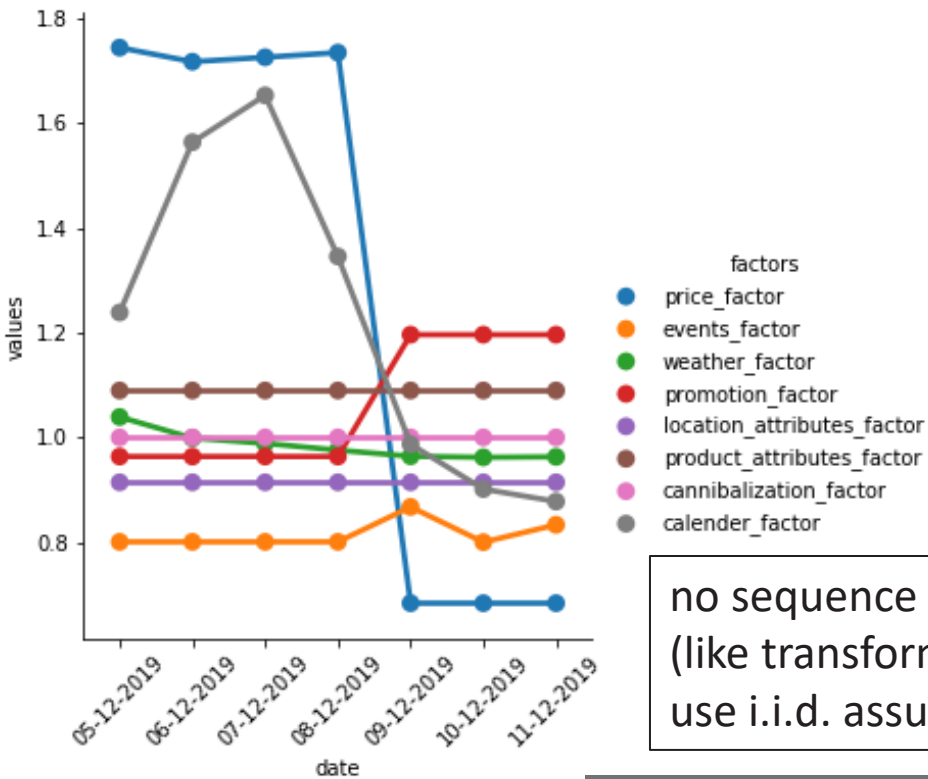- demand (approximately) following Poisson (or rather negative binomial) distribution

# Cyclic Boosting - Prediction View | Individual Explainability

Cyclic Boosting in multiplicative regression mode

**multiplicative model**

variation proportional to level

**individual item-store-day predictions**

Cyclic Boosting allows for detailed explanation of each individual prediction by means of contributions (in form of factors) of each feature in the model.

prediction: look up learned factors of relevant bin for each feature

$$\hat{y}_i = \mu \cdot \prod_{j=1}^{p} f_j^k \quad \text{with} \quad k = \{x_{j,i} \in b_j^k\}$$

global target average

product over factors for all $p$ features in corresponding bins of sample $i$

bin $k$ of feature $j$

factors
- price_factor
- events_factor
- weather_factor
- promotion_factor
- location_attributes_factor
- product_attributes_factor
- cannibalization_factor
- calender_factor

no sequence model (like transformers), use i.i.d. assumption

data binning of features (think of histograms)

do not confuse explainability with causality though
→ need for causal assumptions (e.g., specific smoothing)

# Cyclic Boosting Training
## Coordinate Descent: Boosting-like Update of Factors

```
while (stop criteria)      iterations
...
    for (features)         sequential
    ...
        for (samples)      parallel
```

1. calculate global average μ, initialize all factors to 1

2. cyclically iterate through features and calculate factors for each feature bin (corresponding to minimization of quadratic loss)

multiplicative regression mode (other modes work accordingly)

bin $k$

numerator: target values

factors for corresponding feature bins of sample $i$

for simplicity: show only non-aggregated mode

$$f_j^k = \sum_{x_{j,i} \in b_j^k} y_i \bigg/ \sum_{x_{j,i} \in b_j^k} \hat{y}_{j,i} \quad \bigg| \quad \hat{y}_{j,i} = \mu \cdot \prod_{l \neq j} f_l^k$$

feature $j$

product over all features excluding $j$

sum over all samples $i$ in bin $k$ of feature $j$

denominator: predictions excluding factor from current feature

3. stop according to MAD or MSE criteria at end of iterations (full feature cycles) or when reaching given maximal number of iterations
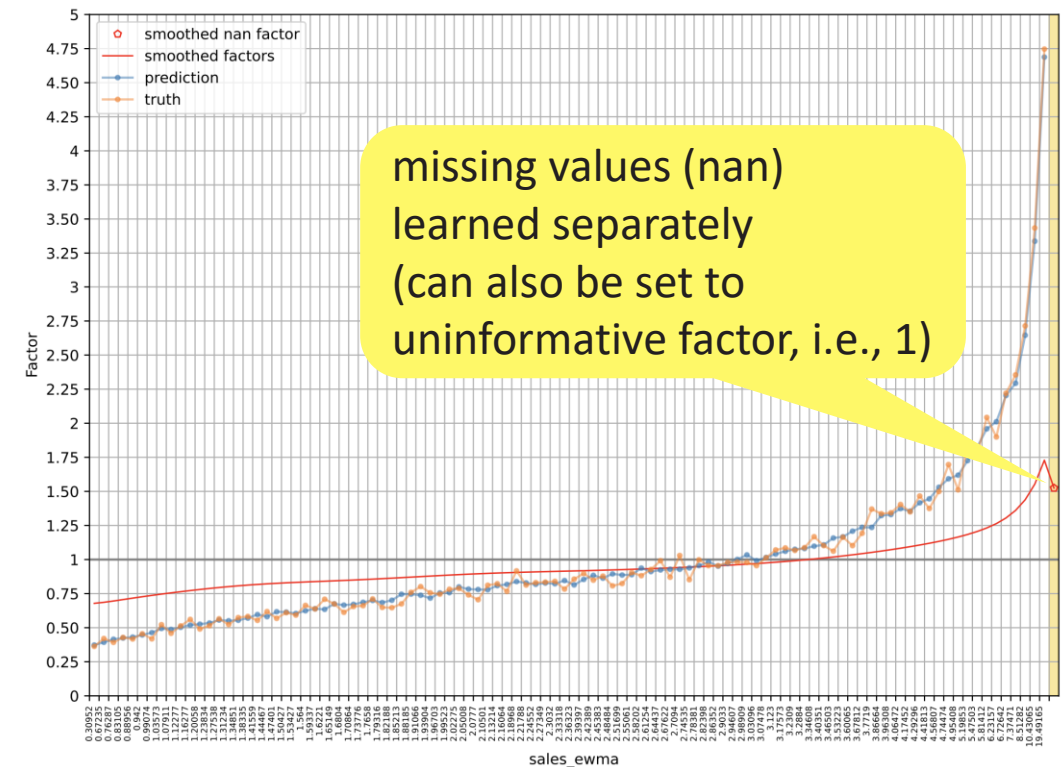
10

# Binning

categorical features retain original categories (learning of specific factor for each of the bins)
→ supporting categorical features with high cardinality

continuous features discretized to:
- either having same bin width (equidistant binning)
- or containing approximately same number of observations (equistatistics binning) with different bin widths



truth = average of target values for all the observed training samples in this bin divided by the global average of target values

y axis: factors for each bin

x axis: e.g., product IDs as bins

missing values (nan) learned separately (can also be set to uninformative factor, i.e., 1)

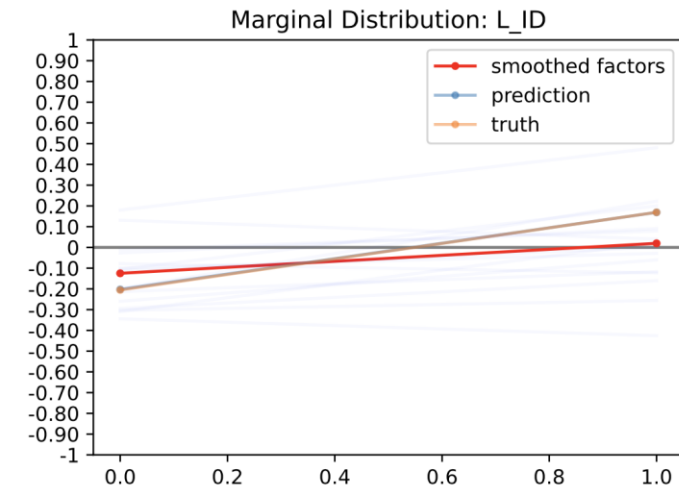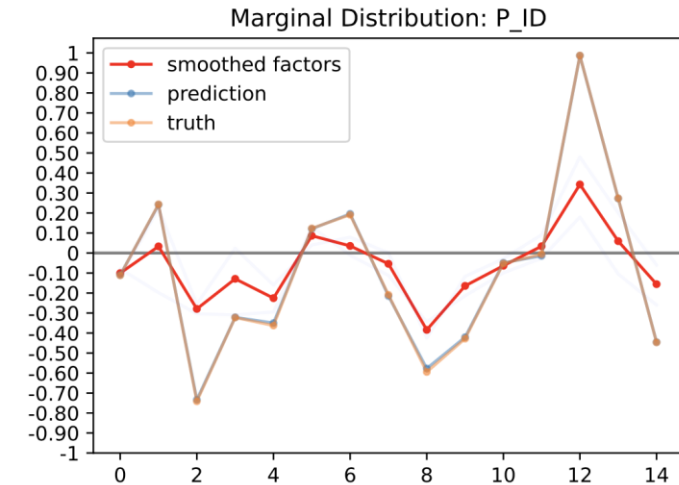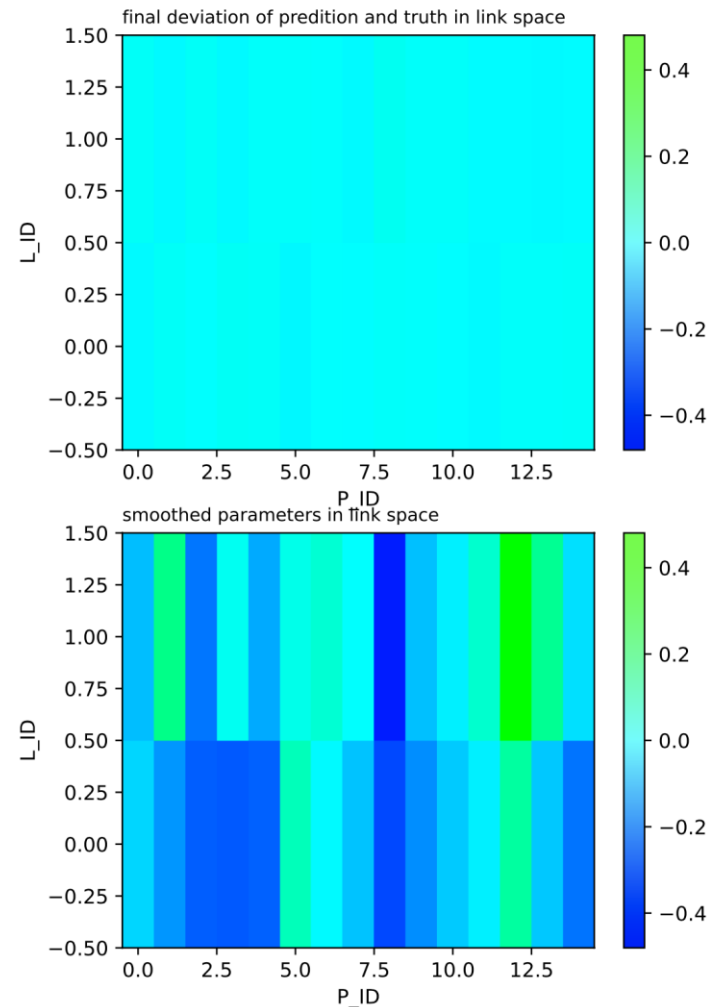local optimization in each bin: allows learning of rare effects with low bias

# Interaction Terms

e.g., different holiday effects for different products

require even more local optimization (rare effects)

→ include binned interaction terms (e.g., 2D or 3D)

can (partly) enable hierarchical model structure (in combination with coordinate descent): interaction terms with product groups, products, locations

# Smoothing

to avoid overfitting:
regularization (smoothing) across bins
→ **drastic reduction of variance** by ignoring fluctuations
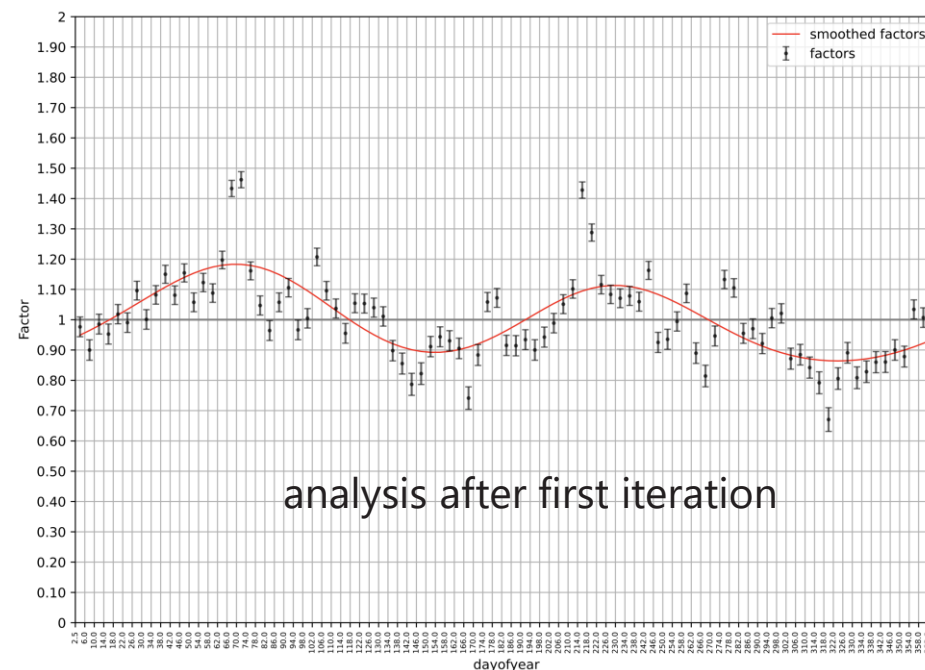
separate smoothings in each iteration
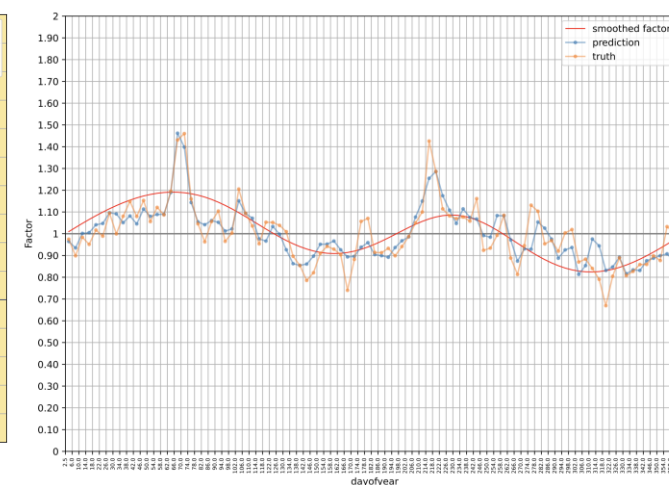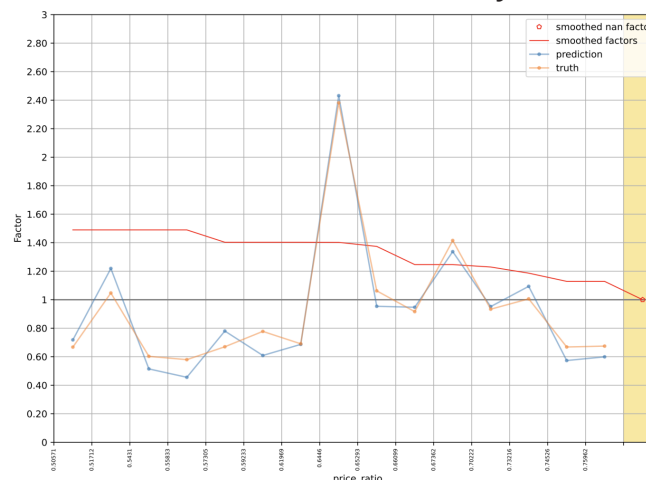
in general: orthogonal polynomials

another way to include *prior knowledge* via
- monotonic requirements
- sinusoidal functions
- (piecewise) linear

use fitted functions (smoothed factors) instead of original factors



analysis after first iteration

analysis after last iteration:

13

# Analysis Plots

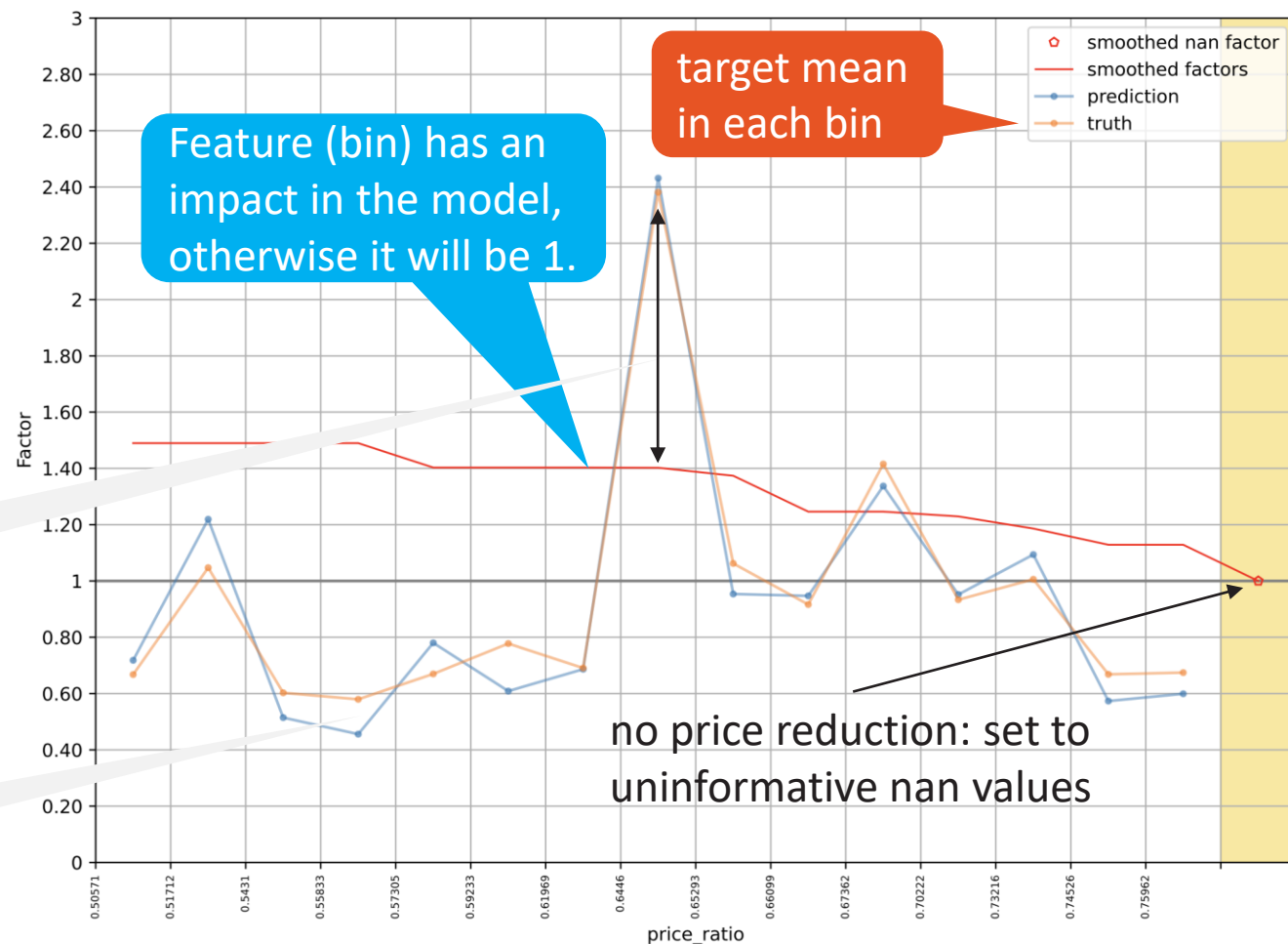**support EDA and modeling**

**model transparency → ease of development**
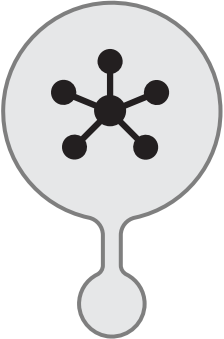
**automatically generated**

Deviations between smoothed factors and predictions come from correlations with other features.

Deviations of predictions from truths show potential model weaknesses (biases) in different bins.
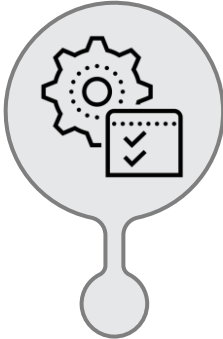
analysis after last iteration:



Feature (bin) has an impact in the model, otherwise it will be 1.

target mean in each bin

no price reduction: set to uninformative nan values

# Cyclic Boosting | Characteristics

| TOPOLOGY | OPTIMIZATION | PREREQUISITE | REGULARIZATION |
|---|---|---|---|
| Generalized Additive Model | cyclic coordinate descent | binning of features | Bayesian updates and smoothing of factor distributions |

similarities to:    backfitting, forward-stagewise modeling (aka boosting)    LightGBM

# Off-The-Shelf Method for Structured Data

- "simple" algorithm, but robust and fast

- few hyperparameters to be tuned

- not much data pre-processing needed

- easily configurable for different data types

- supporting missing values in input data

- assisting model development with individual analysis plots for features

- allowing buiding of complex models by means of interaction terms


- (multiplicative or additive) regression (location parameter)

- classification

# Other Modes: Width Prediction (Scale)

important busines application: automated replenishment

**full, individual PDF predictions (e.g., probability distributions for each product-location-day combination)**

**by means of separate ML models for mean and variance (actually, indirect prediction of variance via dispersion parameter), assuming negative binomial distribution of target (e.g., demand) in maximum likelihood estimation**



### 1. Forecast Probabilities

Product   Sales   Promo   ...

Probability
15%
10%
5%
0%
   0   5   10   15   20   25
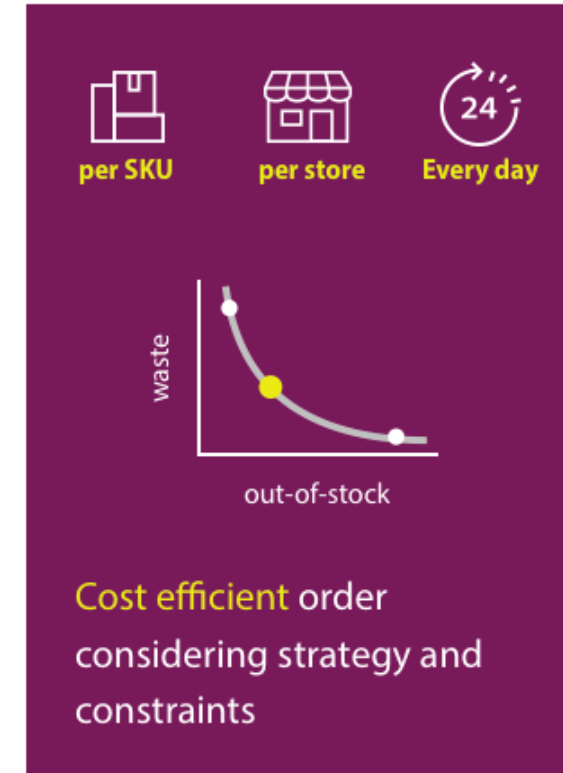Demand (SKU, store, day)

Weather   Holidays   Events   ...

We understand internal & external factors. By forecasting the probability density, we know the risks of e.g. lost sales vs. waste.

### 2. Optimize Decisions

Strategy by cost/benefit

Reduce lost sales   VS.   Reduce shrink   VS.   ...

Constraints
• Case rounding
• Min/max order quantities
• ...

### 3. Automate Orders

per SKU   per store   Every day (24)

waste / out-of-stock

Cost efficient order considering strategy and constraints

Knowing these risks we calculate the order, which minimizes these risks and balances them according to strategy set by the retailer.
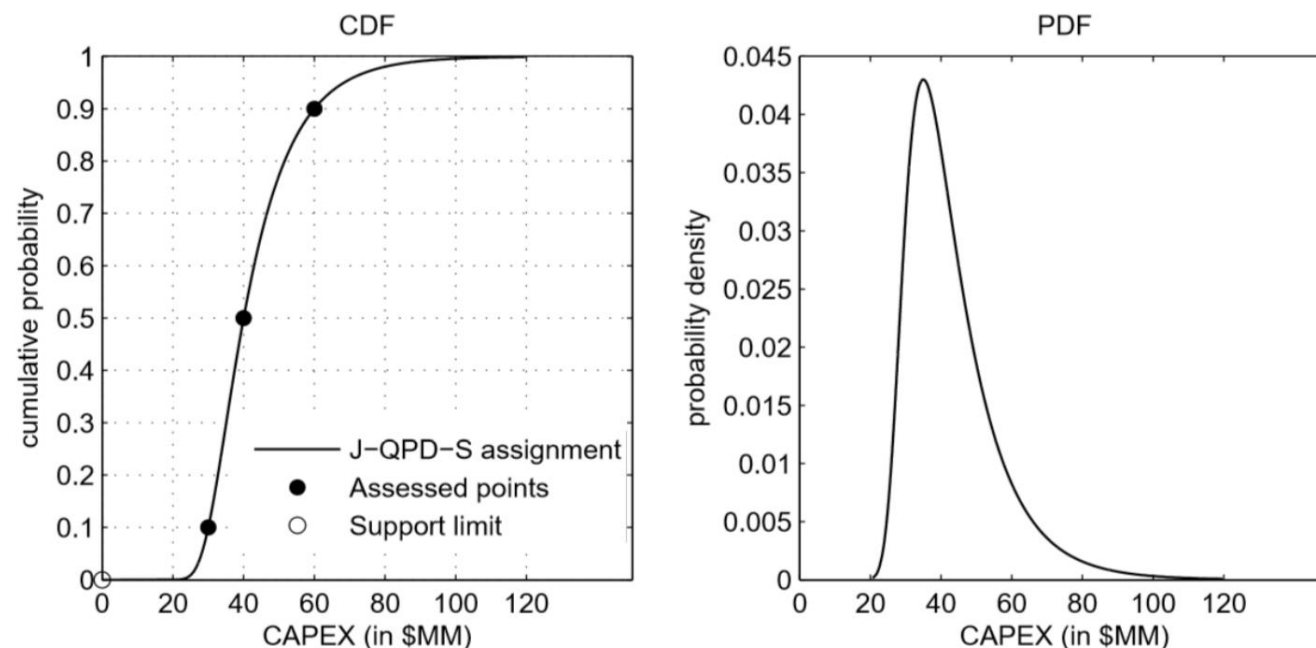
# Alternative: Quantile-Parameterized Distributions (QPD)

**idea: approximate full, individual probability distribution for each sample by using estimated quantiles**

- e.g., from Cyclic Boosting's quantile regression mode (or any other quantile regression method)

- quantiles as parameters of smooth distribution → no fitting, no strict functional assumption

**Johnson QPDs (J-QPD) are parameterized by symmetric quantile triplet**

**implemented in Cyclic Boosting**
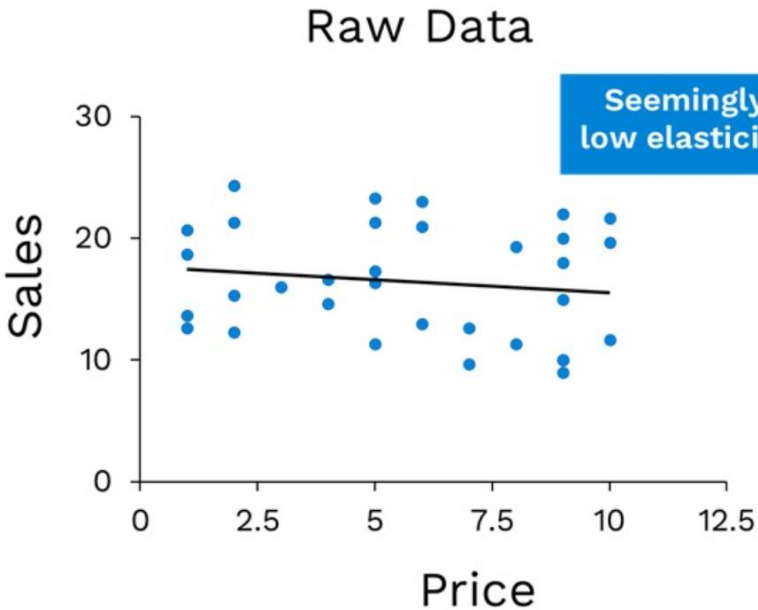
J-QPD

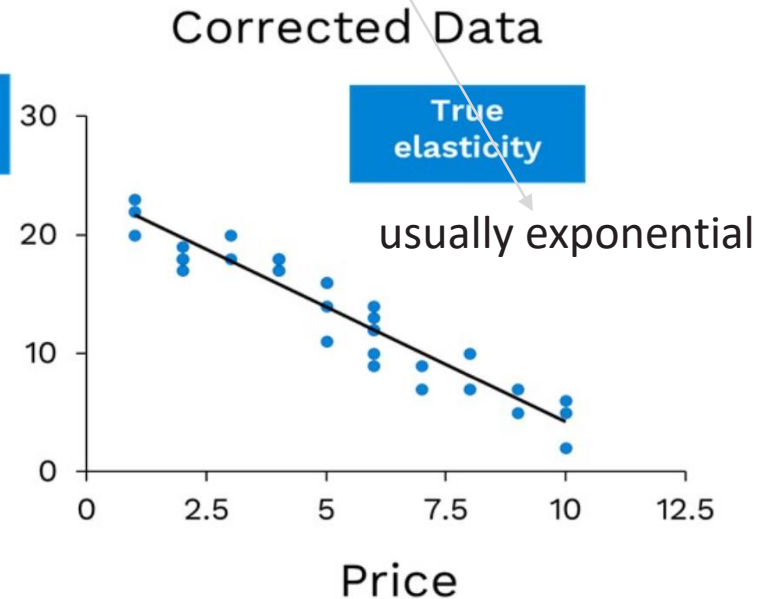# Other Modes: Elasticity Prediction (Shape), Background Subtraction

**other business application: demand shaping by causal inference (mainly beyond ML/CB), examples:**

- **dynamic pricing: influence demand of different products by price setting**

- **customer targeting: influence individual customer demand by couponing (subtract unaffected customers)**

confounded effect:                  after de-confounding:                  use for pricing policies:



**Raw Data**

Seemingly low elasticity

Sales / Price



**Corrected Data**

True elasticity

usually exponential

Price



Maximize items sold

Maximize revenue

Maximize profit

Items Sold / Price