

# Modeling Promotional Pasta Sales with Hierarchical Poisson Autoregression

Group 2

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Background . . . . .	2
1.2	Goals . . . . .	2
<b>2</b>	<b>Data Description</b>	<b>2</b>
<b>3</b>	<b>Model Derivation</b>	<b>2</b>
3.1	Extended Poisson Autoregressive Model . . . . .	2
3.2	Likelihood Function . . . . .	2
3.3	Gradient (Score Function) . . . . .	3
<b>4</b>	<b>Estimation Methods</b>	<b>3</b>
4.1	EM Algorithm . . . . .	3
4.2	Newton-Raphson . . . . .	3
4.3	Bayesian Inference (MCMC) . . . . .	3
<b>5</b>	<b>Model Evaluation and Predictive Performance</b>	<b>4</b>
<b>6</b>	<b>Machine Learning Benchmarks</b>	<b>4</b>
6.1	Random Forest . . . . .	4
6.2	Hidden Markov Model . . . . .	4
<b>7</b>	<b>Summary and Discussion</b>	<b>4</b>

## 1 Introduction

We analyze pasta sales data from an Italian grocery store, consisting of 116 items across four brands, observed daily from 2014 to 2018. Our goal is to model item-level sales with dependence on past observations, promotional status, and latent brand effects.

We compare likelihood-based models (EM, Newton-Raphson), Bayesian inference (MCMC), and machine learning benchmarks (Random Forests, Hidden Markov Models), using predictive accuracy as the main evaluation metric.

## 1.1 Background

## 1.2 Goals

# 2 Data Description

```
data("data_set_tidy")
glimpse(data_set_tidy)
```

```
## Rows: 212,164
## Columns: 5
## $ DATE <date> 2014-01-02, 2014-01-02, 2014-01-02, 2014-01-02, 2014-01-02, 201~
## $ brand <chr> "B1", "B1", "B1", "B1", "B1", "B1", "B1", "B1", "B1", "B1", "B1"~
## $ item <chr> "1", "2", "3", "4", "5", "6", "7", "8", "9", "10", "11", "12", "~
## $ QTY <dbl> 7, 3, 0, 2, 3, 1, 0, 4, 0, 0, 0, 0, 7, 2, 5, 2, 0, 7, 1, 0, 0, 4~
## $ PROMO <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
```

Each row corresponds to daily sales for a specific item, along with brand and promotion indicators.

## 3 Model Derivation

### 3.1 Extended Poisson Autoregressive Model

Let the sales outcome for item  $i$  at time  $t$  be denoted by  $y_{it}$ , where:

$$i \in \{1, \dots, n\}, \quad t \in \{1, \dots, T\}, \quad g_i \in \{1, \dots, B\}$$

Here,  $g_i$  denotes the brand (group) to which item  $i$  belongs.

We define the model as:

$$y_t \sim \text{Poisson}(m_t)$$

$$m_t = \sum_{l=1}^q \beta_l y_{t-l} + \left(1 - \sum_{l=1}^q \beta_l\right) \cdot \exp(\mathbf{x}_t^\top \boldsymbol{\gamma})$$

Where:

- $\beta_l$ : AR coefficients (constrained so sum  $\leq 1$ )
- $\boldsymbol{\gamma}$ : covariate effect vector

### 3.2 Likelihood Function

Let  $\boldsymbol{\theta}$  denote the full set of model parameters. The full likelihood across all items and times is:

$$\mathcal{L}(\boldsymbol{\theta}) = \prod_{t=q+1}^T \frac{m_t^{y_t} e^{-m_t}}{y_t!}$$

The corresponding log-likelihood, used for optimization and posterior inference, is:

$$\log \mathcal{L}(\boldsymbol{\theta}) = \sum_{t=q+1}^T [y_t \log(m_t) - m_t - \log(y_t!)]$$

Where,

$$m_t = \underbrace{\sum_{l=1}^q \beta_l y_{t-l}}_{\text{AR part}} + \underbrace{\left(1 - \sum_{l=1}^q \beta_l\right)}_{\text{mixing weight}} \cdot \underbrace{\exp(\mathbf{x}_t^\top \boldsymbol{\gamma})}_{\text{covariate part}}$$

### 3.3 Gradient (Score Function)

Define the following:

- $a_t = \sum_{l=1}^q \beta_l y_{t-l}$
- $c_t = \exp(\mathbf{x}_t^\top \boldsymbol{\gamma})$
- $w = 1 - \sum \beta_t$
- $m_t = a_t + w \cdot c_t$

Then, the derivative of the log-likelihood with respect to  $\gamma_j$  is:

$$\frac{\partial \log \mathcal{L}}{\partial \gamma_j} = \sum_{t=q+1}^T \left[ \frac{y_t}{m_t} - 1 \right] \cdot w \cdot c_t \cdot x_{tj}$$

And the derivative of the log-likelihood with respect to  $\beta_k$  is:

$$\frac{\partial \log \mathcal{L}}{\partial \beta_k} = \sum_{t=q+1}^T \left[ \frac{y_t}{m_t} - 1 \right] [y_{t-k} - c_t]$$

## 4 Estimation Methods

### 4.1 EM Algorithm

- E-step: Applicable if we add new random intercepts (either item or brand level)
- M-step: Maximize the expected log-likelihood with respect to model parameters  $\gamma$ ,  $\beta$ , and  $\tau$

### 4.2 Newton-Raphson

- Use gradients and Hessians of the log-likelihood to iteratively update parameters

### 4.3 Bayesian Inference (MCMC)

- Fit the independent Poisson Autoregression (PAR) model for each item with MCMC (currently implemented in STAN)
- Place priors directly on model parameters
- $\tau$ : mixing weight for AR component,  $\beta = \tau \tilde{\beta}$

Priors:

- $\boldsymbol{\gamma} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma_\gamma)$
- $\Sigma_\gamma \sim \text{Inv-Wishart}(\nu, \Psi)$  — (approximated as fixed in Stan)
- $\tilde{\beta} \sim \text{Dirichlet}(\alpha)$
- $\tau \sim \text{Beta}(a, b)$

## 5 Model Evaluation and Predictive Performance

We evaluate models using holdout-based prediction: the final 10–20% of each item’s time series is reserved for testing.

```
# Example: results <- evaluate_holdout(...)
# knitr::kable(results)
```

**\*\* Evaluation Metrics \*\*** - Root Mean Squared Error (RMSE) - Mean Absolute Error (MAE) - Log Predictive Density (for Bayesian models) - Posterior predictive interval coverage (for MCMC)

## 6 Machine Learning Benchmarks

### 6.1 Random Forest

```
# rf_fit <- fit_random_forest(...)
```

### 6.2 Hidden Markov Model

```
# hmm_fit <- fit_hmm(...)
```

## 7 Summary and Discussion

- Which methods give the most accurate predictions for sparse, autocorrelated data?
- What is the trade-off between hierarchical shrinkage and item-specific flexibility?
- Are machine learning models (e.g., random forests) better suited for forecasting than structured generative models?