

Competence Development

Time Series Forecasting Review

Jonas Petersen
Ander Runge Walther

Revision History

Revision	Date	Author(s)	Description
0.0.1	21-08-2024	JP	Created

Contents

Chapter 1

Overview

1.1 Introduction

Take a benchmark dataset for forecasting and test out different forecasting algorithms and their strengths and weaknesses. The goal is to get light experience with a broad range of algorithms and/or Python packages and their syntax. The algorithms of interest include

1. Gaussian Processes (via PyMC),
2. Bayesian Two-layer perceptron (via PyMC and tensorflow or keras),
3. Kalman Filter (filterpy?)
4. Decision trees (XGBoost, LightGBM, PyMC)
5. Prohpet,
6. LSTM (tensorflow or keras?),
7. SARIMAX (via statsmodels),
8. Exponential smoothing (via statsmodels)

VS Code with Anaconda will be used. A Git project for each competence development project will be created.

1.2 Data

Chapter 2

Profile Approach

2.1 Method

Consider the problem of having a time series that yield the demand of some product with one entry per week and generating a representative forecast for this. The problem can be framed as a game between Nature and a Robot, where each make a decision. Nature's decision ($s_i \in \Omega_S$) is to pick the true demand and the Robot's decision ($U_i(D) \in \Omega_U$) is guided to minimize a cost function defined viz

$$C : \Omega_U \times \Omega_S \mapsto \mathbb{R} \quad (2.1)$$

where Ω_U denotes the action space of the Robot and Ω_S denotes the action space of Nature. The objective of the Robot is to minimize the expected cost across its actions. Let

$$G \equiv \sum_{i=1}^n C(U_i(D), S_i), \quad (2.2)$$

then

$$\mathbb{E}[G|D, I] = \int ds_1 ds_2 \dots ds_n \sum_{i=1}^n C(U_i(D), s_i) p(s_1, s_2, \dots s_n | D, I) \quad (2.3)$$

where $i \in$ "weeks in forecast horizon" cover the weeks in the forecast horizon and I denotes the background information. The minimum cost across the Robot's actions is defined viz

$$\left. \frac{d}{dU_i(D)} \mathbb{E}[G|D, I] \right|_{U_i(D)=U_i^*(D)} \stackrel{!}{=} 0. \quad (2.4)$$

where $U_i^*(D)$ is the optimal decision rule for the Robot. In this project, the cost function is taken to be

$$C(U_i(D), s_i) = (U_i(D) - s_i)^2 \quad (2.5)$$

Combining Equation 2.3, Equation 2.5 and Equation 2.4 yields

$$U_i^*(D) = \int ds_1 ds_2 \dots ds_n s_i p(s_1, s_2, \dots s_n | D, I). \quad (2.6)$$

Equation 2.6 informs that the optimal decisions for the Robot should be the expected decisions of Nature, which makes intuitive sense considering the cost function (equation (2.5)).

To apply equation Equation 2.6, the probability distribution $p(s_1, s_2, \dots s_n|D, I)$ needs to be specified. To do this, a host of assumptions regarding the process generating the data must be made.

2.1.1 Simple Poisson Assumption

A fixed decision rule that does not require updating at a later point it preferred by the stakeholders.

Axiom 1 (Static and Independent). *The demand (s_i) for a given week of the year belong to the same static, Poisson distribution with a distinct rate parameter for each week.*

Example 1.

According to axiom 1, the sales from each week for a given brand and sub category belong to the same, static distribution. E.g. data from week 5 from 2021, 2022,... belong to a static distribution, and data from week 4 from 2021, 2022,... belong to another static distribution.

Using axiom 1

$$\begin{aligned} p(s_1, s_2, \dots s_n|D, I) &= p(s_1|s_2, \dots s_n, D, I)p(s_2|\dots s_n, D, I) \dots p(s_n|D, I) \\ &= \prod_{i=1}^n p(s_i|D, I), \end{aligned} \quad (2.7)$$

meaning the optimal decision rule for the Robot can be written

$$\begin{aligned} U_i^*(D) &= \mathbb{E}[S_i|D, I] \\ &= \int ds_i s_i p(s_i|D, I) \end{aligned} \quad (2.8)$$

According to axiom 1

$$\begin{aligned} p(s_i|D, I) &= \int d\lambda_i p(s_i, \lambda_i|D, I) \\ &= \int d\lambda_i p(s_i|\lambda_i, D, I)p(\lambda_i|D, I) \\ &= \int d\lambda_i p(s_i|\lambda_i, D, I) \frac{p(D|\lambda_i, I)p(\lambda_i|I)}{p(D|I)} \end{aligned} \quad (2.9)$$

where

$$p(D|I) = \int d\lambda_i p(D|\lambda_i, I)p(\lambda_i|I) \quad (2.10)$$

and

$$\begin{aligned} p(s_i|\lambda_i, D, I) &= p(s_i|\lambda_i, I) \\ &= \text{Poi}(s_i|\lambda_i). \end{aligned} \quad (2.11)$$

Note that the rate parameter λ_i is associated to the demand s_i , whereas the data contains measurements of sales y_j for j representing past weeks. Thus, in order to formulate the likelihood function, a relationship between sales y and demand s must be assumed. In general

Definition 1 (Demand and Sales). *Let $m_{b,c,i}$ denote the stock of a given brand b and sub category c for week i , then*

$$y_{b,c,i} = \begin{cases} s_{b,c,i} & \text{if } s_{b,c,i} \leq m_{b,c,i} \\ m_{b,c,i} & \text{else} \end{cases}. \quad (2.12)$$

Axiom 2 (Never Out of Stock). *For simplicity, products are assumed to be never of of stock, meaning (referring to definition 1) $y_{b,c,i} = s_{b,c,i}$.*

Using axioms 1 and 2, the likelihood function can be written

$$p(D|\lambda_i, I) = \prod_{j \in \text{week } i} \text{Poi}(y_j|\lambda_i), \quad (2.13)$$

where only data from the corresponding week is used in accordance with axiom 1. An obvious choice of prior would be the conjugate Gamma distribution (which is also the distribution with maximum entropy) with parameters α, β , meaning

$$p(\lambda_i|I) = \text{Ga}(\lambda_i|\alpha, \beta) \quad (2.14)$$

This means

$$\frac{p(D|\lambda_i, I)p(\lambda_i|I)}{p(D|I)} = \text{Ga}(\lambda_i|\alpha + N\bar{y}, \beta + n) \quad (2.15)$$

where

$$\bar{y} \equiv \frac{1}{n} \sum_{j \in \text{week } i} y_j. \quad (2.16)$$

Then

$$\begin{aligned} \mathbb{E}[s_i|D, I] &= \sum_{s_i} s_i \int d\lambda_i \text{Poi}(s_i|\lambda_i) \text{Ga}(\lambda_i|\alpha + N\bar{y}, \beta + n_i) \\ &= \int d\lambda_i \lambda_i \text{Ga}(\lambda_i|\alpha + N\bar{y}, \beta + n_i) \\ &= \frac{\alpha + n_i\bar{y}_i}{\beta + n_i} \end{aligned} \quad (2.17)$$

In the limit of $\alpha, \beta \rightarrow 0$

$$\begin{aligned} U_i^*(D) &= \lim_{\alpha, \beta \rightarrow 0} \mathbb{E}[S_i|D, I] \\ &= \bar{y}_i, \end{aligned} \quad (2.18)$$

Thus, using axioms 1 and 2, a fixed decision rule yielding the forecasted demand for week i in the future, represented by Equation 2.18, is obtained.

2.1.2 Advanced Poisson Assumption

Axiom 3 (Static with coupled rate parameter). *The demand (s_i) for a given week of the year belong to the same static Poisson Distribution where the rate parameter is coupled between different weeks.*

Using axiom 3

$$\begin{aligned} p(s_1, s_2, \dots, s_n|D, I) &= \int d\theta p(s_1, s_2, \dots, s_n, \theta|D, I) \\ &= \int d\theta p(s_1, s_2, \dots, s_n|\theta, D, I)p(\theta|D, I) \end{aligned} \quad (2.19)$$

where

$$\begin{aligned} p(s_1, s_2, \dots, s_n|\theta, D, I) &= p(s_1, s_2, \dots, s_n|\theta, I) \\ &= \prod_{j \in \text{forecast horizon}} \text{Poi}(s_j|\lambda(\theta, j)) \end{aligned} \quad (2.20)$$

and

$$p(\theta|D, I) = \frac{p(D|\theta, I)p(\theta|I)}{p(D|I)} \quad (2.21)$$

$$p(D|\theta, I) = \prod_{i \in \text{training data}} \text{Poi}(s_i|\lambda(\theta, i)). \quad (2.22)$$

This means

$$\begin{aligned} U_j^*(D) &= \int d\theta \int ds_1 ds_2 \dots ds_n s_j \text{Poi}(s_j|\lambda(\theta, j))p(\theta|D, I) \\ &= \int d\theta \int ds_j s_j \text{Poi}(s_j|\lambda(\theta, j))p(\theta|D, I) \\ &= \int d\theta \lambda(\theta, j)p(\theta|D, I) \\ &= \mathbb{E}[\lambda|D, I] \end{aligned} \quad (2.23)$$

where λ is some model, in this case a Fourier series

$$\lambda(\theta, j) = a_0 + \sum_{l=1}^M \left(a_l \cos \frac{2\pi}{N} j + b_l \sin \frac{2\pi}{N} j \right) \quad (2.24)$$

and $\theta = \{a, b\}$.

Chapter 3

PYMC Approach

The advantage is that there is complete control over the architecture and the formalism is Bayesian, meaning there is uncertainty and a statistical model behind. The downside is the computation time and the expertise required to get a good model running.

3.1 pymc TLP

The advantage is this is a flexible one-size-fits-all approach, which can accommodate multidimensional time series. The downside is that the computation time and lack of interpretability of the coefficients. Feature engineering is also required to get a good result.

3.2 pymc prophet

The advantage is a flexible model that can accommodate multidimensional time series and no feature engineering is required. The downside is the computation time and difficulty in getting the result to converge. The coefficients are highly interpretable, which is a strength, but this leans into the difficulty of obtaining convergence of a good result.

3.3 pymc ar

3.4 pymc gaussian processes

Chapter 4

Statsmodels

4.1 TimeSeriesAnalysis

The advantage is it does not require feature engineering, is fairly accurate and fast. It also includes uncertainty.

Chapter 5

Light GBM Approach

This is very fast and accurate, however, with not apparent uncertainty quantification and it does not generalize to multiple time series.