

## THE M5 COMPETITION

### Competitors' Guide

#### Contents

<b>Objectives</b> .....	2
<b>Dates and hosting</b> .....	3
<b>The dataset</b> .....	5
<b>Evaluation</b> .....	8
Forecasting horizon.....	8
Point forecasts .....	8
Probabilistic forecasts .....	10
Weighting.....	12
<b>The Prizes</b> .....	15
Distribution of prize money .....	15
Reproducibility .....	15
<b>Publications</b> .....	16
<b>The Benchmarks</b> .....	16
Point forecasts .....	16
Probabilistic forecasts .....	20
<b>Submission</b> .....	21

# M5

## Objectives

The objective of the M5 forecasting competition is to advance the theory and practice of forecasting by identifying the method(s) that provide the **most accurate point forecasts** for each of the **42,840** time series of the competition. In addition, to elicit information to estimate the **uncertainty distribution** of the realized values of these series as precisely as possible.

### 목표

M5 예측 경쟁의 목표는 42,840 개의 각 경쟁 시계열에 대해 가장 정확한 포인트 예측을 제공하는 방법을 식별하여 예측 이론 및 실습을 발전시키는 것입니다. 또한 이 시리즈의 실현 된 값의 불확실성 분포를 가능한 정확하게 추정하기 위한 정보를 이끌어 내기 위해.

To that end, the participants of M5 are asked to provide **28 days ahead point forecasts (PFs)** for all the series of the competition, as well as the corresponding **median and 50%, 67%, 95%, and 99% prediction intervals (PIs)**.

이를 위해 M5 참가자에게는 모든 경쟁 시리즈에 대해 28 일의 사전 예측 (PF)과 해당 중간 값 및 50 %, 67 %, 95 % 및 99 % 예측 간격을 제공해야 합니다 (PI).

The M5 differs from the previous four ones in five important ways, some of them suggested by the discussants of the M4<sup>1</sup> competition, as follows:

- First, it uses **grouped** unit sales data, starting at the product-store level and being aggregated to that of product departments, product categories, stores, and three geographical areas: the States of California (CA), Texas (TX), and Wisconsin (WI).
- Second, besides the time series data, it includes **explanatory variables** such as sell prices, promotions, days of the week, and special events (e.g. Super Bowl, Valentine's Day, and Orthodox Easter) that typically affect unit sales and could improve forecasting accuracy.
- Third, in addition to point forecasts, it assesses the **distribution of uncertainty**, as the participants are asked to provide information on nine indicative quantiles.
- Fourth, instead of having a single competition to estimate both the point forecasts and the uncertainty distribution, there will be **two** parallel tracks using the **same** dataset, the first requiring 28 days ahead point forecasts and the second 28 days ahead probabilistic forecasts for the median and four prediction intervals (50%, 67%, 95%, and 99%).
- Fifth, for the first time it focuses on series that display **intermittency**, i.e., sporadic demand including zeros.

---

<sup>1</sup> Makridakis, S., Spiliotis, E. & Assimakopoulos, V. (2020). The M4 Competition: 100,000 time series and 61 forecasting methods. International Journal of Forecasting, 36, 54–74.

# M5

M5 는 다섯 가지 중요한 점에서 이전 네 가지와 다르며, 그중 일부는 다음과 같이 M4 경쟁 토론자가 제안했습니다.

- 먼저, 제품 매장 수준에서 시작하여 제품 부서, 제품 카테고리, 상점 및 3 개의 지리적 영역 (캘리포니아 주, 캘리포니아 주, 텍사스 주, 위스콘신 (WI)).
- 둘째, 시계열 데이터 외에도 판매 가격, 프로모션, 요일 및 단위 판매에 영향을 미치고 예측 정확도를 향상시킬 수 있는 특별 이벤트 (예 : Super Bowl, Valentine 's Day 및 Orthodox Easter)와 같은 설명 변수가 포함됩니다..
- 셋째, 포인트 예측 외에도 참가자들이 9 개의 지표에 대한 정보를 제공하도록 요청함에 따라 불확실성의 분포를 평가합니다.
- 넷째, 포인트 예측과 불확실성 분포를 추정하기위한 단일 경쟁 대신 동일한 데이터 세트를 사용하는 두 개의 병렬 트랙이 있습니다. 및 4 개의 예측 구간 (50 %, 67 %, 95 % 및 99 %).
- 다섯째, 처음으로 간헐성, 즉 0 을 포함한 산발적 수요를 표시하는 시리츠에 중점을 둡니다.

## Dates and hosting

The M5 will start on **March 2, 2020** and finish on **June 30** of the same year. The competition will be run using the **Kaggle** platform. Thus, we expect many submissions from all types of forecasters including data scientists, statisticians, and practitioners, expanding the field of forecasting and eventually integrating its various approaches for improving accuracy and uncertainty estimation.

### 날짜와 호스팅

M5 는 2020 년 3 월 2 일에 시작하여 같은 해 6 월 30 일에 끝납니다. 경쟁은 Kaggle 플랫폼을 사용하여 진행됩니다. 따라서 우리는 데이터 과학자, 통계 학자 및 실무자를 포함한 모든 유형의 예측자가 제출하여 예측 범위를 확장하고 궁극적으로 정확성 및 불확실성 추정을 개선하기위한 다양한 접근 방식을 통합 할 것으로 기대합니다.

The competition will be divided into two separate Kaggle competitions, using the same dataset, with the first (**M5 Forecasting Competition – Accuracy**) requiring 28 days ahead point forecasts and the second (**M5 Forecasting Competition – Uncertainty**) 28 days ahead probabilistic forecasts for the corresponding median and four prediction intervals (50%, 67%, 95%, and 99%).

경쟁은 동일한 데이터 세트를 사용하여 두 개의 개별 Kaggle 경쟁으로 구분되며, 첫 번째 (M5 예측 경쟁 – 정확성)는 28 일 사전 예측이 필요하고 두 번째 (M5 예측 경쟁 – 불확실성)는 해당 예측에 대한 확률 예측이 28 일입니다. 중앙값 및 4 개의 예측 간격 (50 %, 67 %, 95 % 및 99 %).

# M5

In order to support the participants to validate their forecasting approaches, the competition will include a **validation phase** that will take place from **March 2, 2020 to 31 May** of the same year. During this phase, the participants will be allowed to train their forecasting methods with the data initially provided by the organizers and validate the performance of their approaches using a hidden sample of 28 days, not made publicly available. By submitting their forecasts at the Kaggle platform (a maximum of 5 entries per day), the participants will be informed about the score of their submission, which will be then published in Kaggle's real-time leaderboard. Given this instant feedback, participants will be allowed to effectively revise and resubmit their forecasts by learning from the received feedback.

참가자들에게 그들의 예측 접근법의 유효성을 검증 할 수 있도록하기 위해 경쟁에는 2020 년 3 월 2 일부터 같은 해 5 월 31 일까지 진행되는 검증 단계가 포함됩니다. 이 단계에서 참가자는 주최자가 처음 제공 한 데이터로 예측 방법을 교육하고 공개되지 않은 28 일의 숨겨진 샘플을 사용하여 접근 방식의 성능을 검증 할 수 있습니다. Kaggle 플랫폼 (하루에 최대 5 개 항목)에 예측을 제출하면 참가자에게 제출 점수에 대한 정보가 제공되며 Kaggle 의 실시간 리더 보드에 게시됩니다. 이러한 즉각적인 피드백이 주어지면 참가자는 수신 된 피드백을 통해 예측을 효과적으로 수정하고 다시 제출할 수 있습니다.

After the end of the validation phase, i.e., from **June 1, 2020 to 30 June** of the same year, the participants will be provided with the actual values of the 28 days of data used for scoring their performance during the validation phase. They will be asked then to re-estimate or adjust (if needed) their forecasting models in order to submit their final forecasts and prediction intervals for the following 28 days, i.e., the data used for the final evaluation of the participants. During this time, there will be no leaderboard, meaning that no feedback will be given to the participants about their score after submitting their forecasts. Thus, although the participants will be free to (re)submit their forecasts any time they wish (a maximum of 5 entries per day), they will not be aware of their absolute, as well as their relative performance. The **final ranks** of the participants will be made available only **at the end of competition**, when the test data will be made available. This is done in order for the competition to simulate reality as closely as possible, given that in real life forecasters do not know the future.

같은 해 2020 년 6 월 1 일부터 6 월 30 일까지 유효성 검사 단계가 끝난 후 참가자에게는 유효성 검사 단계에서 성과를 평가하는 데 사용 된 28 일간의 데이터의 실제 값이 제공됩니다. 그런 다음 28 일 동안의 최종 예측 및 예측 간격, 즉 참가자의 최종 평가에 사용 된 데이터를 제출하기 위해 예측 모델을 다시 추정하거나 조정 (필요한 경우)하도록 요청합니다. 이 기간 동안 리더 보드가 없으므로 예측을 제출 한 후 참가자에게 점수에 대한 피드백이 제공되지 않습니다. 따라서 참가자는 원하는 시간 (하루에 최대 5 개의 항목)을 언제든지 자유롭게 예측을 다시 제출할 수 있지만 상대적인 성과뿐만 아니라 절대적인 성과도 알 수 없습니다. 시험 데이터가 제공 될 때, 참가자의 최종 순위는 경쟁이 끝날 때만 이용 가능합니다. 현실에서 예측자가 미래를 알지 못한다는 점을 감안할 때 경쟁이 현실을 최대한 가깝게 시뮬레이션하기 위해 수행됩니다.

Note that the submission system will be open from the beginning of the competition, meaning that participants will be able to submit their final forecast from March 2, 2020 to June 30, 2020, even during the validation phase. However, as previously mentioned, the complete M5 training sample (including the 28 days used for the validation phases' leaderboard) will only become available on June 1, 2020. So, any

# M5

participant submitting his/his/their final forecasts during the validation phase will be missing the last 28 days of the complete training sample.

제출 시스템은 경쟁이 시작될 때부터 시작됩니다. 즉, 참가자는 유효성 검사 단계에서도 2020 년 3 월 2 일부터 2020 년 6 월 30 일까지 최종 예측을 제출할 수 있습니다. 그러나 앞에서 언급 한 바와 같이 완전한 M5 교육 샘플 (검증 단계의 리더 보드에 사용 된 28 일 포함)은 2020 년 6 월 1 일에 만 사용할 수 있습니다. 따라서 모든 참가자는 유효성 검사 단계에서 최종 예측을 제출합니다. 전체 교육 샘플의 지난 28 일이 누락됩니다.

Note also that M5 will be divided into **two tracks**, one requiring point forecasts, and one requiring the estimation of the uncertainty distribution, each with its separate prizes of \$50,000. Thus, two individual competitions will be visible at the Kaggle platform, each one with its own separate leaderboard. Participants are allowed to compete and be eligible for prizes in the first track, the second track, or both.

M5 는 두 개의 트랙으로 나뉘며, 하나는 포인트 예측이 필요하고, 하나는 불확실성 분포의 추정이 필요하며, 각각은 별도의 상금이 \$ 50,000 입니다. 따라서 Kaggle 플랫폼에서 두 개의 개별 경기가 표시되며, 각 개별 경기는 별도의 리더 보드가 있습니다. 참가자는 첫 번째 트랙, 두 번째 트랙 또는 둘 다에서 상을받을 수 있습니다.

## The dataset

The M5 dataset, generously made available by **Walmart**, involves the unit sales of various products sold in the USA, organized in the form of **grouped time series**. More specifically, the dataset involves the unit sales of **3,049 products**, classified in **3 product categories** (Hobbies, Foods, and Household) and **7 product departments**, in which the above-mentioned categories are disaggregated. The products are sold across **ten stores**, located in **three States** (CA, TX, and WI). In this respect, the bottom-level of the hierarchy, i.e., product-store unit sales can be mapped across either product categories or geographical regions, as follows:

**Table 1: Number of M5 series per aggregation level.**

Level id	Aggregation Level	Number of series
1	Unit sales of all products, aggregated for all stores/states	1
2	Unit sales of all products, aggregated for each State	3
3	Unit sales of all products, aggregated for each store	10
4	Unit sales of all products, aggregated for each category	3
5	Unit sales of all products, aggregated for each department	7
6	Unit sales of all products, aggregated for each State and category	9
7	Unit sales of all products, aggregated for each State and department	21
8	Unit sales of all products, aggregated for each store and category	30
9	Unit sales of all products, aggregated for each store and department	70
10	Unit sales of product x, aggregated for all stores/states	3,049
11	Unit sales of product x, aggregated for each State	9,147

12	Unit sales of product x, aggregated for each store	30,490
Total		42,840

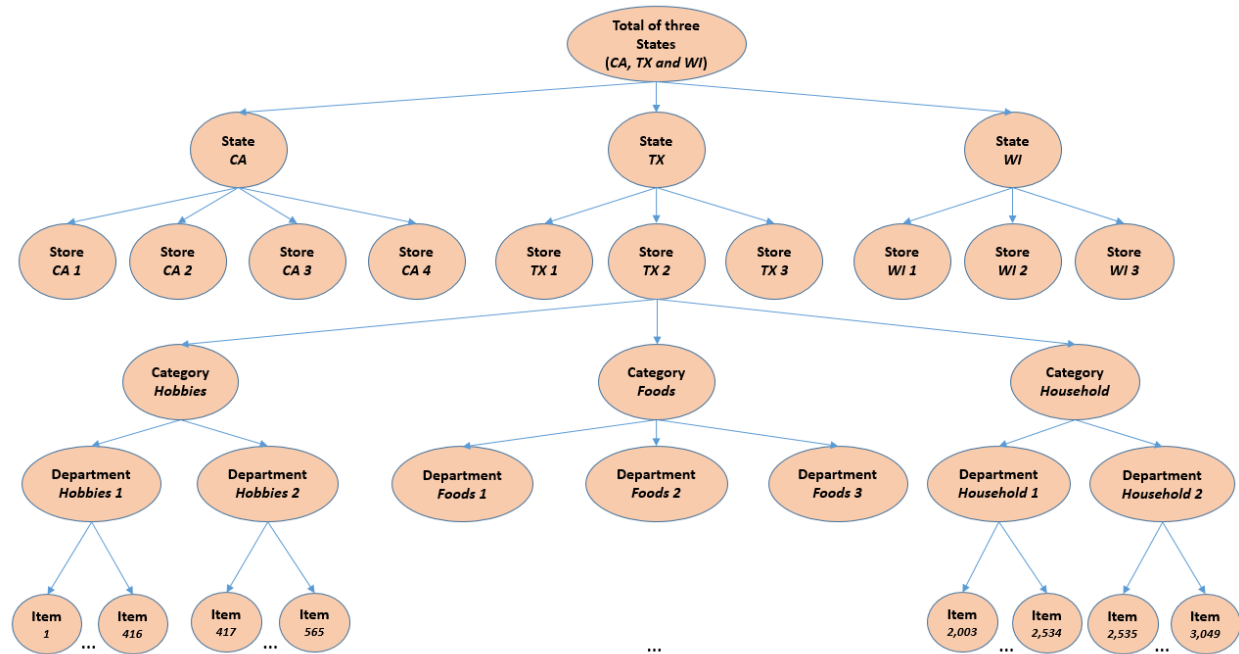


Figure 1: An overview of how the M5 series are organized.

Walmart 가 관대하게 제공하는 M5 데이터 세트에는 미국에서 판매 된 다양한 제품의 단위 판매가 그룹화 된 시계열 형식으로 구성되어 있습니다. 보다 구체적으로, 데이터 집합에는 3 가지 제품 범주 (취미, 음식 및 가구)로 분류 된 3,049 개의 제품과 7 개의 제품 부서 (위에서 언급 한 범주가 분리됨)의 단위 판매가 포함됩니다. 이 제품은 3 개 주 (CA, TX 및 WI)에있는 10 개의 매장에서 판매됩니다. 이와 관련하여 다음과 같이 계층의 최하위 레벨, 즉 제품 상점 단위 판매를 제품 카테고리 또는 지역에 매핑 할 수 있습니다.

The historical data range from **2011-01-29** to **2016-06-19**. Thus, the products have a (maximum) selling history of 1,941<sup>2</sup> days / 5.4 years (**test data of h=28 days not included**).

The M5 dataset consists of the following **three (3) files**:

<sup>2</sup> This number refers to the history of the final training dataset to be provided at June 1, 2020. 28 less days will be available during the validation phase, as explained in the “dates and hosting” section.

# M5

과거 데이터는 2011 년 1 월 29 일부터 2016 년 6 월 19 일까지입니다. 따라서, 제품은 1,941 일 / 5.4 년의 (최대) 판매 이력을 가지고 있습니다 ( $h = 28$  일의 테스트 데이터는 포함되지 않음).

M5 데이터 세트는 다음 3 개의 파일로 구성됩니다.

## File 1: “calendar.csv”

Contains information about the dates the products are sold.

- *date*: The date in a “y-m-d” format.
- *wm\_yr\_wk*: The id of the week the date belongs to.
- *weekday*: The type of the day (Saturday, Sunday, ..., Friday).
- *wday*: The id of the weekday, starting from Saturday.
- *month*: The month of the date.
- *year*: The year of the date.
- *event\_name\_1*: If the date includes an event, the name of this event.
- *event\_type\_1*: If the date includes an event, the type of this event.
- *event\_name\_2*: If the date includes a second event, the name of this event.
- *event\_type\_2*: If the date includes a second event, the type of this event.
- *snap\_CA*, *snap\_TX*, and *snap\_WI*: A binary variable (0 or 1) indicating whether the stores of CA, TX or WI allow SNAP<sup>3</sup> purchases on the examined date. 1 indicates that SNAP purchases are allowed.

## File 2: “sell\_prices.csv”

Contains information about the price of the products sold per store and date.

- *store\_id*: The id of the store where the product is sold.
- *item\_id*: The id of the product.
- *wm\_yr\_wk*: The id of the week.
- *sell\_price*: The price of the product for the given week/store. The price is provided per week (average across seven days). If not available, this means that the product was not sold during the examined week. Note that although prices are constant at weekly basis, they may change through time (both training and test set).

## File 3: “sales\_train.csv”

Contains the historical daily unit sales data per product and store.

- *item\_id*: The id of the product.
- *dept\_id*: The id of the department the product belongs to.
- *cat\_id*: The id of the category the product belongs to.
- *store\_id*: The id of the store where the product is sold.
- *state\_id*: The State where the store is located.

---

<sup>3</sup> The United States federal government provides a nutrition assistance benefit called the Supplement Nutrition Assistance Program (SNAP). SNAP provides low income families and individuals with an Electronic Benefits Transfer debit card to purchase food products. In many states, the monetary benefits are dispersed to people across 10 days of the month and on each of these days 1/10 of the people will receive the benefit on their card. More information about the SNAP program can be found here: <https://www.fns.usda.gov/snap/supplemental-nutrition-assistance-program>

# M5

- $d_1, d_2, \dots, d_i, \dots, d_{1941}$ : The number of units sold at day  $i$ , starting from 2011-01-29.

## Evaluation

### Forecasting horizon

The number of forecasts required, both for point and probabilistic forecasts, is  **$h=28$  days** (4 weeks ahead).

The performance measures are **first computed for each series** separately by averaging their values across the forecasting horizon and **then averaged again across the series** in a weighted fashion (see below) to obtain the final scores.

### 평가

### 예측 지평

포인트 예측과 확률 예측에 필요한 예측 수는  $h = 28$  일 (4 주 전)입니다.

성능 측정은 먼저 예측 기간에 걸쳐 값을 평균하여 각 계열에 대해 개별적으로 계산된 다음 가중치를 기준으로 계열 전체에 대해 평균을 다시 평균하여 (아래 참조) 최종 점수를 얻습니다.

### Point forecasts

The accuracy of the point forecasts will be evaluated using the **Root Mean Squared Scaled Error (RMSSE)**, which is a variant of the well-known Mean Absolute Scaled Error (MASE) proposed by Hyndman and Koehler (2006)<sup>4</sup>. The measure is calculated for each series as follows:

### 포인트 예측

점 예측의 정확도는 Hyndman 과 Koehler (2006)가 제안한 잘 알려진 평균 절대 스케일 오차 (MASE)의 변형 인 RMSSE (루트 평균 제곱 스케일 오차)를 사용하여 평가됩니다. 각 계열에 대한 측정 값은 다음과 같이 계산됩니다.

$$RMSSE = \sqrt{\frac{1}{h} \frac{\sum_{t=n+1}^{n+h} (Y_t - \hat{Y}_t)^2}{\frac{1}{n-1} \sum_{t=2}^n (Y_t - Y_{t-1})^2}},$$

where  $Y_t$  is the actual future value of the examined time series at point  $t$ ,  $\hat{Y}_t$  the generated forecast,  $n$  the length of the training sample (number of historical observations), and  $h$  the forecasting horizon.

<sup>4</sup> R. J. Hyndman & A. B. Koehler (2006). Another look at measures of forecast accuracy. International Journal of Forecasting 22(4), 679-688.



# M5

여기서  $y_t$  는 시점  $t$  에서 검사 된 시계열의 실제 미래 값,  $(y_t)$  generated 생성 된 예측,  $n$  훈련 샘플의 길이 (이력 관측치 수) 및  $h$  예측 수평선입니다.

Note that the denominator of RMSSE is computed only for the time-periods for which the examined product(s) are actively sold, i.e., the periods following the first non-zero demand observed for the series under evaluation.

RMSSE 의 분모는 검사 된 제품이 활발하게 판매 된 기간, 즉 평가중인 시리즈에 대해 0 이 아닌 첫 번째 수요 이후의 기간에 대해서만 계산됩니다.

The choice of the measure is justified as follows:

- The M5 series are characterized by intermittency, involving sporadic unit sales with lots of zeros. This means that absolute errors, which are optimized for the median, would assign lower scores (better performance) to forecasting methods that derive forecasts close to zero. However, the objective of M5 is to accurately forecast the average demand and for this reason, the accuracy measure used builds on squared errors, which are optimized for the mean.
- The measure is scale independent, meaning that it can be effectively used to compare forecasts across series with different scales.
- In contrast to other measures, it can be safely computed as it does not rely on divisions with values that could be equal or close to zero (e.g. as done in percentage errors when  $Y_t = 0$  or relative errors when the error of the benchmark used for scaling is zero).
- The measure penalizes positive and negative forecast errors, as well as large and small forecasts, equally, thus being symmetric.

측정의 선택은 다음과 같이 정당화됩니다.

M5 시리즈는 간헐적으로 특징이 있으며 간헐적으로 단위 판매가 0 이됩니다. 이는 중앙값에 최적화된 절대 오차가 0 에 가까운 예측을 도출하는 예측 방법에 낮은 점수 (성능 향상)를 할당 함을 의미합니다. 그러나 M5 의 목표는 평균 수요를 정확하게 예측하는 데 사용되므로 정확도 측정 값은 평균에 최적화 된 제품 오차를 기반으로합니다.

이 척도는 척도와 무관하므로 여러 척도를 기준으로 계열 간 예측을 효과적으로 비교하는 데 효과적으로 사용할 수 있습니다.

다른 측정 값과 달리, 0 과 같거나 0 에 가까운 값을 갖는 나누기에 의존하지 않으므로 안전하게 계산할 수 있습니다 (예 :  $y_t = 0$  일 때 백분율 오류로 수행되거나 벤치 마크 오류가 사용 된 경우 상대 오류로 수행됨) 스케일링은 0 입니다).

이 측정 값은 크고 작은 예측뿐만 아니라 양수 및 음수 예측 오류에 동일하게 적용되므로 대칭입니다.

# M5

After estimating the RMSSE for all the 42,840 time series of the competition, the participating methods will be ranked using the **Weighted RMSSE (WRMSSE)**, as described latter in this Guide, using the following formula:

42,840 번의 모든 경쟁 시계열에 대해 RMSSE 를 추정 한 후 참여 방법은 다음 공식을 사용하여 가이드의 뒷부분에 설명 된대로 WRMSSE (Weighted RMSSE)를 사용하여 순위가 매겨집니다.

$$WRMSSE = \sum_{i=1}^{42,840} w_i * RMSSE,$$

where  $w_i$  is the weight of the  $i_{th}$  series of the competition. A lower WRMSSE score is better.

Note that the weight of each series will be computed based on the last 28 observations of the training sample of the dataset, i.e., the cumulative actual dollar sales that each series displayed in that particular period (sum of units sold multiplied by their respective price). An indicative example for computing the WRMSSE will be available on the GitHub<sup>5</sup> repository of the competition.

여기서  $w_i$  는  $i_{th}$  시리즈 경쟁의 가중치입니다. WRMSSE 점수가 낮을수록 좋습니다.

각 시리즈의 가중치는 데이터 세트의 훈련 샘플에 대한 마지막 28 개의 관측치, 즉 특정 기간에 각 시리즈가 표시 한 누적 실제 달러 판매량 (판매 된 단위의 합에 각각의 가격을 곱한 값)을 기반으로 계산됩니다. . WRMSSE 를 계산하는 대표적인 예는 경쟁사의 GitHub 저장소에서 사용할 수 있습니다.

## Probabilistic forecasts

The precision of the probabilistic forecasts will be evaluated using the **Scaled Pinball Loss (SPL)** function. The measure is calculated for each series and quantile as follows:

### 확률 적 예측

확률 예측의 정확도는 SPL (Scaled Pinball Loss) 기능을 사용하여 평가됩니다. 측정은 각 시리즈에 대해 계산되고 다음과 같이 Quantile 됩니다

$$SPL(u) = \frac{\frac{1}{h} \sum_{t=n+1}^{n+h} (Y_t - Q_t(u)) u \mathbf{1}\{Q_t(u) \leq Y_t\} + (Q_t(u) - Y_t)(1 - u) \mathbf{1}\{Q_t(u) > Y_t\}}{\frac{1}{n-1} \sum_{t=2}^n |Y_t - Y_{t-1}|},$$

where  $Y_t$  is the actual future value of the examined time series at point  $t$ ,  $Q_t(u)$  the generated forecast for quantile  $u$ ,  $h$  the forecasting horizon,  $n$  the length of the training sample (number of historical observations), and  $\mathbf{1}$  the indicator function (being 1 if  $Y$  is within the postulated interval and 0 otherwise).

As done with RMSSE, the denominator of SPL is computed only for the time-periods for which the examined items/products are actively sold, i.e., the periods following the first non-zero demand observed for the series under evaluation.

<sup>5</sup> <https://github.com/Mcompetitions>

# M5

여기서  $y_t$  는 점  $t$  에서 검사 된 시계열의 실제 미래 값이며,  $Q_t(u)$   $u$  에 대한 생성 된 예측  $u$ , 예측 수평선,  $n$  훈련 샘플의 길이 (이력 관측치 수) 및 1 함수 ( $y$  가 가정 된 간격 내에 있으면 1 이되고 그렇지 않으면 0 이 됨).

RMSSE 와 마찬가지로 SPL 의 분모는 검사 된 품목 / 제품이 활발하게 판매되는 기간, 즉 평가중인 시리즈에 대해 첫 번째 0 이 아닌 수요 다음 기간에 대해서만 계산됩니다.

Given that forecasters will be asked to provide the **median**, and the **50%, 67%, 95%, and 99% PIs**,  $u$  is set to  $u_1=0.005$ ,  $u_2=0.025$ ,  $u_3=0.165$ ,  $u_4=0.25$ ,  $u_5=0.5$ ,  $u_6=0.75$ ,  $u_7=0.835$ ,  $u_8=0.975$ , and  $u_9=0.995$ . The smaller values of  $u$  correspond to the left side of the distribution, while the higher values to the right side of the distribution, with  $u = 0.5$  being the median. The median and the 50% and 67% PIs provide a good sense of the middle of the distribution, while the 95% and 99% PIs provide information about its tails, which are important in terms of the risk of extremely high or extremely low outcomes.

After estimating the SPL for all the 42,840 time series of the competition and for all the requested quantiles, the participating methods will be ranked using the **Weighted SPL (WSPL)**, as described latter in this Guide, divided by nine (average performance of nine quantiles across all series) , using the following formula:

예측 자에게 중간 값과 50 %, 67 %, 95 % 및 99 % PI 를 제공하도록 요청하면  $u$  는  $u_1 = 0.005$ ,  $u_2 = 0.025$ ,  $u_3 = 0.165$ ,  $u_4 = 0.25$ ,  $u_5 = 0.5$  로 설정됩니다. ,  $u_6 = 0.75$ ,  $u_7 = 0.835$ ,  $u_8 = 0.975$  및  $u_9 = 0.995$  입니다.  $u$  의 작은 값은 분포의 왼쪽에 해당하고,  $u$  의 값은 분포의 오른쪽에,  $u = 0.5$  는 중앙값입니다. 중앙값 및 50 % 및 67 % PI 는 분포 중간에 대한 좋은 의미를 제공하는 반면 95 % 및 99 % PI 는 꼬리에 대한 정보를 제공하며 결과는 매우 높거나 매우 낮은 결과의 위험 측면에서 중요합니다.

42,840 번의 모든 경쟁 시계열 및 요청 된 모든 Quantile 에 대해 SPL 을 추정 한 후, 본 가이드의 후반부에 설명 된대로 가중 SPL (WSPL)을 사용하여 참여 방법의 순위를 9 로 나눕니다 (평균 9 개의 Quantile). 다음 공식을 사용하여 모든 시리즈에서) :

$$WSPL = \sum_{i=1}^{42,840} w_i * \frac{1}{9} \sum_{j=1}^9 SPL(u_j),$$

where  $w_i$  is the weight of the  $i_{th}$  series of the competition and  $u_j$  the  $j_{th}$  out of the examined quantiles. A lower WSPL score is better.

여기서  $w_i$  는  $i_{th}$  시리즈 경쟁의 가중치이고  $u_j$  는 검사 된 Quantile 의  $j_{th}$  out 입니다. WSPL 점수가 낮을수록 좋습니다.

The choice of the measure is justified as follows:

# M5

- PL is scaled in a similar fashion to that of RMSSE, meaning that it can be effectively used to compare forecasts across series with different scales. Moreover, SPL can be safely computed as it does not rely on divisions with values that could be equal to zero.
- Since M5 does not focus on a particular decision-making problem, neither defines the exact parameters of such a problem (which could also vary for different aggregation levels and series), it becomes evident that all quantiles could be potentially useful. Moreover, since the objective of the M5 is to estimate the uncertainty distribution of the realized values of the examined series as precisely as possible, both sides and both ends of the distribution are considered relevant. In this regard, no special weights are assigned to the examined quantiles, which are therefore equally weighted.

측정의 선택은 다음과 같이 정당화됩니다.

- PL 은 RMSSE 와 유사한 방식으로 조정됩니다. 즉, 일련의 예측을 다른 척도와 효과적으로 비교하는데 효과적으로 사용할 수 있습니다. 또한 SPL 은 0 과 같은 값을 갖는 나누기에 의존하지 않으므로 안전하게 계산할 수 있습니다.
- M5 는 특정 의사 결정 문제에 중점을 두지 않기 때문에 이러한 문제의 정확한 매개 변수를 정의하지 않으며 (다른 집계 수준 및 시리즈에 따라 다를 수도 있음) 모든 Quantile 이 잠재적으로 유용 할 수 있음이 분명해집니다. 더욱이, M5 의 목적은 가능한 한 정확하게 검사 된 계열의 실현 된 값의 불확실성 분포를 추정하는 것이기 때문에 분포의 양측과 양단이 적절한 것으로 간주됩니다. 이와 관련하여, 검사 된 Quantile 에 특별한 가중치가 할당되지 않으므로 동일한 가중치가 적용됩니다.

Note that, once again, the weight of each series will be computed based on the last 28 observations of the training sample of the dataset, i.e., the cumulative actual dollar sales that each series displayed in that particular period (sum of units sold multiplied by their respective price). An indicative example for computing the WSPL will be available on the GitHub repository of the competition.

다시 한 번, 각 시리즈의 가중치는 데이터 세트의 훈련 샘플에 대한 마지막 28 개의 관측치, 즉 특정 기간에 각 시리즈가 표시 한 누적 실제 달러 판매량 (판매 된 판매량의 합)을 기반으로 계산됩니다. 해당 가격). WSPL 을 계산하는 대표적인 예는 경쟁사의 GitHub 리포지토리에서 사용할 수 있습니다.

## Weighting

In contrast to the previous M competition, M5 involves the unit sales of various products of different selling volumes and prices that are organized in a hierarchical fashion. This means that, businesswise, in order for a method to perform well, it must provide accurate forecasts across all hierarchical levels, especially for series of high importance, i.e. for series that represent significant sales, measured in US dollars. In other words, we expect from the best performing forecasting methods to derive lower forecasting errors for the series that are more valuable for the company.

가중치

# M5

이전 M 경쟁과 달리 M5 는 판매량과 가격이 다른 다양한 제품의 판매 단위를 계층 구조로 구성합니다. 즉, 비즈니스 측면에서, 방법이 잘 수행 되려면 모든 계층 적 수준에서, 특히 일련의 높은 중요도, 즉 상당한 매출을 나타내는 일련의 미국 달러로 정확한 예측을 제공해야 합니다. 다시 말해, 우리는 최고 성능의 예측 방법을 통해 회사에 더 유용한 시리즈의 예측 오류를 줄일 수 있을 것으로 기대합니다.

To that end, the forecasting errors computed for each participating method (both RMSSE and SPL) will be weighted across the M5 series based on their **cumulative actual dollar sales**, which is a good and objective proxy of their actual value for the company in monetary terms. The cumulative dollar sales will be computed using **the last 28 observations of the training sample** (sum of units sold multiplied by their respective price), i.e., a period equal to the forecasting horizon. Note that since both the number of units being sold and their respective price change through time, this estimation is based on the sum of the corresponding daily dollar sales.

이를 위해 각 참여 방법 (RMSSE 및 SPL 모두)에 대해 계산 된 예측 오류는 누적 실제 달러 판매를 기준으로 M5 시리즈에 가중치가 적용됩니다. . 누적 달러 판매는 교육 샘플의 마지막 28 개 관측치 (판매 된 단위의 합계에 각 가격을 곱한 값), 즉 예측 기간과 동일한 기간을 사용하여 계산됩니다. 판매되는 단위 수와 각각의 가격이 시간에 따라 변하므로 이 추정치는 해당 일일 달러 판매의 합계를 기반으로 합니다.

Below you may find a simple, yet indicative example of how these weights will be computed:  
아래에서 이러한 가중치를 계산하는 방법에 대한 간단하지만 설명적인 예를 찾을 수 있습니다.

Assume that two products of the same department, A and B, are sold in a store at WI and we are interested in forecasting the unit sales of these two products, as well as their aggregate sales. Thus, in this example, we consider two different aggregation levels ( $K=2$ ), the first level consisting of two series (series A and B) and the second level of a single series (sum of series A and B).

동일한 부서의 두 제품 A 와 B 가 WI 의 상점에서 판매되고 있으며 두 제품의 판매량과 총 판매량을 예측하는 데 관심이 있다고 가정하십시오. 따라서 이 예에서는 두 개의 서로 다른 집계 수준 ( $K = 2$ )을 고려합니다. 첫 번째 수준은 두 계열 (시리즈 A 및 B)과 두 번째 수준의 단일 계열 (시리즈 A 및 B 의 합)입니다.

Product A displayed a total of \$10 in sales in the last 28 days of the training sample, while product B \$12. Thus, the aggregate dollar sales of products A and B in the last 28 days were \$22. Assume also that a forecasting method was used to derive point forecasts for product A, product B, and their aggregate unit sales, displaying errors  $RMSSE_A=0.8$ ,  $RMSSE_B=0.7$ , and  $RMSSE=0.77$ , respectively. If the M5 dataset involved just those three series, the final WRMSSE score of the method would be

제품 A 는 교육 샘플의 마지막 28 일 동안 총 10 달러의 매출을 보인 반면 제품 B 는 12 달러였습니다. 따라서 지난 28 일 동안 제품 A 와 B 의 총 달러 판매량은 22 달러였습니다. 또한 예측 방법을 사용하여 제품 A, 제품 B 및 총 단위 판매량에 대한 포인트 예측을 도출하고 각각  $RMSSE_A = 0.8$ ,  $RMSSE_B = 0.7$

# M5

및  $RMSSE = 0.77$  을 표시한다고 가정합니다. M5 데이터 세트에이 세 시리즈 만 포함 된 경우 분석법의 최종 WRMSSE 점수는 다음과 같습니다.

$$WRMSSE = RMSSE_A * w_1 + RMSSE_B * w_2 + RMSSE * w_3 =$$

$$RMSSE_A * \frac{1}{K} * \frac{\$ Sales_A}{\$ Sales_A + \$ Sales_B} + RMSSE_B * \frac{1}{K} * \frac{\$ Sales_B}{\$ Sales_A + \$ Sales_B} + RMSSE * \frac{1}{K}$$

$$* \frac{\$ Sales}{\$ Sales_A + \$ Sales_B} =$$

$$0.8 * \frac{1}{2} * \frac{10}{10+12} + 0.7 * \frac{1}{2} * \frac{12}{10+12} + 0.77 * \frac{1}{2} * 1 = 0.758.$$

This weighting scheme can be expanded in order to consider more stores, geographical regions, product categories, and product departments, as previously described. Since the M5 competition involves twelve aggregation levels, K is set equal to 12, with the weights of the series being computed so that they sum to one at each aggregation level.

이 가중치 체계는 앞에서 설명한대로 더 많은 상점, 지역, 제품 범주 및 제품 부서를 고려하기 위해 확장 될 수 있습니다. M5 경쟁에는 12 개의 집계 레벨이 포함되므로 K 는 12 로 설정되며 시리즈의 가중치는 각 집계 레벨에서 하나씩 합산되도록 계산됩니다.

Respectively, RMSSE, which is used in the above equation for estimating WRMSSE, can be replaced with SPL to compute WSPL.

각각 WRMSSE 를 추정하기 위해 위의 방정식에서 사용되는 RMSSE 는 SPL 로 대체되어 WSPL 을 계산할 수 있습니다.

Note that **all hierarchical levels are equally weighted**. The reason is because the total dollar sales of a product, measured across all three States, are equal to the sum of the dollar sales of this product when measured across all ten stores. Similarly, because the total dollar sales of a product category of a store are equal to the sum of the dollar sales of the departments that this category consists of, as well as the sum of the dollar sales of the products of the corresponding departments. Moreover, as previously discussed for the case of the probabilistic forecasts, M5 does not focus on a particular decision-making problem, which means that there is no reason for weighting unequally the individual levels of the hierarchy.

모든 계층 수준에 동일한 가중치가 적용됩니다. 그 이유는 3 개 주 전체에서 측정 된 제품의 총 달러 판매가 10 개 매장 모두에서 측정 할 때이 제품의 달러 판매 합계와 같기 때문입니다. 마찬가지로, 상점의 제품 카테고리에 대한 총 달러 판매는이 카테고리 구성된 부서의 달러 판매의 합계와 해당 부서의 제품에 대한 달러 판매의 합계와 같습니다. 더욱이, 확률 론적 예측의 경우에 대해 이전에 논의 된 바와 같이, M5 는 특정 의사 결정 문제에 초점을 맞추지 않으므로, 계층의 개별 레벨을 동일하게 가중 할 이유가 없음을 의미한다.

# M5

An indicative example for computing WRMSSE and WSPL will be available on the GitHub repository of the competition, indicating among others the exact weight of each series in the competition.

WRMSSE 및 WSPL 을 계산하는 대표적인 예는 경쟁사의 GitHub 저장소에서 사용할 수 있으며, 이는 경쟁에서 각 시리즈의 정확한 무게를 나타냅니다.

## The Prizes

### Distribution of prize money

There will be **12** major prizes awarded to the winners of the M5 Competition, which will be further distributed among the participants based on (i) the hierarchical levels that their forecasts exceeded and (ii) the quantiles of the uncertainty distribution that were better captured. The Prizes will be awarded on **December 8, 2020**, during the **M5 Forecasting Conference** to be held in New York City. At this date, Kaggle will be issuing the payments digitally using its collaborating firm Payoneer.

The total of the \$100,000 prize money will be distributed equally between the Forecasting and Uncertainty M5 competition as follows:

Prize id	Prize	Amount
1A	Most accurate point forecasts	\$25,000
2A	Second most accurate point forecasts	\$10,000
3A	Third most accurate point forecasts	\$5,000
4A	Fourth most accurate point forecasts	\$3,000
5A	Fifth most accurate point forecasts	\$2,000
6A	Most accurate student point forecasts	\$5,000
	<b>Total: M5 Forecasting Competition - Point Forecasts</b>	<b>\$50,000</b>
1B	Most precise estimation of the uncertainty distribution	\$25,000
2B	Second most precise estimation of the uncertainty distribution	\$10,000
3B	Third most precise estimation of the uncertainty distribution	\$5,000
4B	Fourth most precise estimation of the uncertainty distribution	\$3,000
5B	Fifth most precise estimation of the uncertainty distribution	\$2,000
6B	Most precise student estimation of the uncertainty distribution	\$5,000
	<b>Total: M5 Forecasting Competition - Uncertainty Distribution</b>	<b>\$50,000</b>
	<b>Total: M5 Competition</b>	<b>\$100,000</b>

### Reproducibility

The prerequisite for winning any prize will be that the code used for generating the forecasts, with the exception of companies providing forecasting services and those claiming proprietary software, will be put on GitHub, **not later than 14 days after the end of the competition** (i.e., the 14<sup>th</sup> of July, 2020). In addition, there must be instructions on how to exactly reproduce the M5 submitted forecasts. In this regard, individuals and companies will be able to use the code and the instructions provided, crediting the person/group that has developed them, to improve their organizational forecasts.

Companies providing forecasting services and those claiming proprietary software will have to provide the organizers with a detailed description of how their forecasts were made and a source, or execution



# M5

file for reproducing their forecasts. Given the critical importance of objectivity and replicability, such description and file will be mandatory for winning any prize of the competition. An execution file can be submitted in case that the source program needs to be kept confidential, or, alternatively, a source program with a termination date for running it.

After receiving the code/program/files for reproducing the submitted forecasts, the organizers will evaluate their results in terms of reproducibility. Since some methods may involve random initializations, any method that displays a replicability rate higher than 98% will be considered as fully replicable and be awarded the prize, exactly as done in M4. Otherwise, the prize will be given to the next best-performing and fully reproducible submission.

## Publications

Similar to the M3 and M4 competitions, there will be a special issue of the **International Journal of Forecasting (IJF)** exclusively devoted to all aspects of the M5 Competition with special emphasis on what we have learned and how we can use such learning to improve the theory and practice of forecasting as well as expand its usefulness and applicability.

## The Benchmarks

Like done in the M4 competition, there will be benchmark methods, twenty-four (24) for point forecasts, and six (6) for probabilistic ones. As these methods are well known, readily available, and straightforward to apply, the accuracy of the new ones submitted to the M5 Competition must provide superior accuracy in order to be considered and used in practice (taking also into account the computational time it would be required to utilize a more accurate method versus the benchmarks whose computational requirements are minimal).

### Point forecasts

#### Statistical Benchmarks

**1. Naive:** A random walk model, defined as

$$\hat{Y}_{n+i} = Y_n, i = 1, 2, \dots, h.$$

The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**2. Seasonal Naive (sNaive):** Like Naive, but this time the forecasts of the model are equal to the last known observation of the same period in order for it to capture possible weekly seasonal variations. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**3. Simple Exponential Smoothing<sup>6</sup> (SES):** The simplest exponential smoothing model, aimed at predicting series without a trend, defined as

$$\hat{Y}_t = aY_t + (1 - a)\hat{Y}_{t-1}.$$

---

<sup>6</sup> Gardner Jr., E. S. (1985). Exponential smoothing: The state of the art. *Journal of Forecasting*, 4, 1–28.



The smoothing parameter  $\alpha$  is selected from the range [0.1, 0.3] by minimizing the insample mean squared error (MSE) of the model, while the first observation of the series is used for initialization. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**4. Moving Averages (MA):** Forecasts are computed by averaging the last  $k$  observations of the series, as follows

$$\hat{Y}_t = \frac{\sum_{i=1}^k Y_{t-i}}{k},$$

where  $k$  is selected from the range [2, 5] by minimizing the insample MSE. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**5. Croston's method<sup>7</sup> (CRO):** The method proposed by Croston to forecast series that display intermittent demand. The method decomposes the original series into the non-zero demand size  $z_t$  and the inter-demand intervals  $p_t$ , deriving forecasts as follows:

$$\hat{Y}_t = \frac{\hat{z}_t}{\hat{p}_t},$$

where both  $z_t$  and  $p_t$  are predicted using SES. The smoothing parameter of both components is set equal to 0.1. The first observation of the components are used for initialization. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**6. Optimized Croston's method (optCro):** Like CRO, but this time the smoothing parameter is selected from the range [0.1, 0.3], like done with SES, in order to allow for more flexibility. The non-zero demand size and the inter-demand intervals are smoothed separately using (potentially) different  $\alpha$  parameters. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**7. Syntetos-Boylan Approximation<sup>8</sup> (SBA):** A variant of the Croston's method that utilizes a debiasing factor as follows:

$$\hat{Y}_t = 0.95 \frac{\hat{z}_t}{\hat{p}_t}$$

The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

<sup>7</sup> Croston, J. D. (1972). Forecasting and stock control for intermittent demands. Journal of the Operational Research Society, 23, 289–303.

<sup>8</sup> Syntetos, A. A. & Boylan, J. E. (2005). The accuracy of intermittent demand estimates. International Journal of Forecasting, 21, 303–314.

**8. Teunter-Syntetos-Babai method<sup>9</sup> (TSB):** A modification to Croston's method that replaces the inter-demand intervals component with the demand probability,  $d_t$ , being 1 if demand occurs at time  $t$  and 0 otherwise. Similarly to Croston's method,  $d_t$  is forecasted using SES. The smoothing parameters of  $d_t$  and  $z_t$  may differ, exactly as optCRO. The forecast is given as follows:

$$\hat{Y}_t = \hat{d}_t \hat{z}_t$$

The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**9. Aggregate-Disaggregate Intermittent Demand Approach<sup>10</sup> (ADIDA):** Temporal aggregation is used for reducing the presence of zero observations, thus mitigating the undesirable effect of the variance observed in the intervals. ADIDA uses equally sized time buckets to perform non-overlapping temporal aggregation and predict the demand over a pre-specified lead-time. The time bucket is set equal to the mean inter-demand interval. SES is used to obtain the forecasts. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**10. Intermittent Multiple Aggregation Prediction Algorithm<sup>11</sup> (iMAPA):** Another way for implementing temporal aggregation in demand forecasting. However, in contrast to ADIDA that considers a single aggregation level, iMAPA considers multiple ones, aiming at capturing different dynamics of the data. Thus, iMAPA proceeds by averaging the derived point forecasts, generated using SES. The maximum aggregation level is set equal to the maximum inter-demand interval. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**11. Exponential Smoothing<sup>12</sup> - Top-Down (ES\_td):** An algorithm is used to select the most appropriate exponential smoothing model for predicting the top-level series of the hierarchy (level 1 of Table 1), indicated through information criteria. The top-down method is used for reconciliation (based on historical proportions, estimated for the last 28 days).

**12. Exponential Smoothing – Bottom-Up (ES\_bu):** An algorithm is used to select the most appropriate exponential smoothing model for predicting the bottom-level series of the hierarchy (level 12 of Table 1), indicated through information criteria. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**13. Exponential Smoothing with eXplanatory variables (ESX):** Similar to ES, but this time two explanatory variables are used as regressors to improve forecasting accuracy by providing additional information about the future. The first variable is discrete and takes values 0, 1, 2 or 3, based on the number of States

<sup>9</sup> Teunter, R. H., Syntetos, A. A. & Babai, M. Z. (2011). Intermittent demand: Linking forecasting to inventory obsolescence. *European Journal of Operational Research*, 214, 606–615.

<sup>10</sup> Nikolopoulos, K., Syntetos, A. A., Boylan, J. E., Petropoulos, F. & Assimakopoulos, V. (2011). An aggregate-disaggregate intermittent demand approach (ADIDA) to forecasting: an empirical proposition and analysis. *Journal of the Operational Research Society*, 62, 544–554.

<sup>11</sup> Petropoulos, F. & Kourentzes, N. (2015). Forecast combinations for intermittent demand. *Journal of the Operational Research Society*, 66, 914–924

<sup>12</sup> Hyndman, R.J., Koehler, A.B., Snyder, R.D. & Grose, S. (2002). A state space framework for automatic forecasting using exponential smoothing methods. *International Journal of Forecasting*, 18 (3), 439–454.

that allow SNAP purchases on the examined date. The second variable is binary and indicates whether the examined date includes a special event (1) or not (0). The top-down method is used for reconciliation (based on historical proportions, estimated for the last 28 days).

**14. AutoRegressive Integrated Moving Average<sup>13</sup> - Top-Down (ARIMA\_td):** An algorithm is used to select the most appropriate ARIMA model for predicting the top-level series of the hierarchy (level 1 of Table 1), indicated through information criteria. The top-down method is used for reconciliation (based on historical proportions, estimated for the last 28 days).

**15. AutoRegressive Integrated Moving Average – Bottom-Up (ARIMA\_bu):** An algorithm is used to select the most appropriate ARIMA model for predicting the bottom-level series of the hierarchy (level 12 of Table 1), indicated through information criteria. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**16. AutoRegressive Integrated Moving Average with eXplanatory variables (ARIMAX):** Similar to ARIMA, but this time two external variables are used as regressors to improve forecasting accuracy by providing additional information about the future, exactly as done for the case of ESX. The top-down method is used for reconciliation (based on historical proportions, estimated for the last 28 days).

## Machine Learning Benchmarks

**17. Multi-Layer Perceptron (MLP):** A single hidden layer NN of 14 input nodes (last two weeks of available data), 28 hidden nodes, and one output node. The Scaled Conjugate Gradient method is used for estimating the weights that are initialized randomly, while the maximum iterations are set equal to 500. The activation functions of the hidden and output layers are the logistic and linear one, respectively. In total, 10 MLPs are trained to forecast each series and then the median operator is used to average the individual forecasts in order to mitigate possible variations due to poor weight initializations. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**18. Random Forest (RF):** This is a combination of multiple regression trees, each one depending on the values of a random vector sampled independently and with the same distribution. Given that RF averages the predictions of multiple trees, it is more robust to noise and less likely to over-fit on the training data. We consider a total of 500 non-pruned trees and four randomly sampled variables at each split. Bootstrap sampling is done with replacement. Like done in MLP, the last 14 observations of the series are considered for training the model. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**19. Global Multi-Layer Perceptron (GMLP):** Like MLP, but this time, instead of training multiple models, one for each series, a single model that learns across all series is constructed. This is done given that M4 indicated the beneficial effect of cross learning. The last 14 observations of each series are used as inputs, along with information about the coefficient of variation of non-zero demands ( $CV^2$ ) and the average number of time-periods between two successive non-zero demands (ADI). This additional information is used in order to facilitate learning across series of different characteristics. The forecasting method is used

<sup>13</sup> Hyndman, R. & Khandakar Y. (2008). Automatic time series forecasting: the forecast package for R. Journal of Statistical Software, 26, 1-22.

# M5

for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

**20. Global Random Forest (GRF):** Like GMLP, but instead of using an MLP for obtaining the forecasts, a RF is exploited instead. The forecasting method is used for predicting the series of the lowest level of the hierarchy (level 12 of Table 1) and the bottom-up method is then used for reconciliation.

## Combination Benchmarks

**21. Average of ES and ARIMA, as computed using the bottom-up approach (Com\_b):** The simple arithmetic mean of ES\_bu and ARIMA\_bu.

**22. Average of ES and ARIMA, as computed using the top-down approach (Com\_t):** The simple arithmetic mean of ES\_td and ARIMA\_td.

**23. Average of the two ES methods, the first computed using the top-down approach and the second using the bottom-up approach (Com\_tb):** The simple arithmetic mean of ES\_td and ES\_bu.

**24. Average of the global and local MLPs (Com\_lg):** The simple arithmetic mean of MLP and GMLP. The bottom-up method is then used for reconciliation.

Observe that the benchmark methods {1-10, 12, 15, 17-20} are applied at the product-store level of the hierarchically structured dataset. Thus, the bottom-up method is used for obtaining reconciled forecasts for the rest of the hierarchical levels. On the other hand, the benchmark methods {11, 13, 14, 16} are applied at the top level of the hierarchically structured dataset. Thus, the top-down method is used for obtaining reconciled forecasts for the rest of the hierarchical levels (based on historical proportions, estimated for the last 28 days).

## Probabilistic forecasts

**i. Naive:** Similar implementation to the Naive 1 used for computing point forecasts.

**ii. Seasonal Naive (sNaive):** Similar implementation to the sNaive one used for computing point forecasts.

**iii. Simple Exponential Smoothing (SES):** Similar implementation to the SES one used for computing point forecasts.

**iv. Exponential Smoothing (ES):** Similar implementation to the ES\_bu one used for computing point forecasts.

**v. AutoRegressive Integrated Moving Average (ARIMA):** Similar implementation to the ARIMA\_bu one used for computing point forecasts.

**vi. Kernel density estimate (Kernel):** A kernel is used to estimate the corresponding quantiles in the historical data that are then used as probabilistic forecasts.

The code for generating the forecasts of the abovementioned benchmarks will be available on the GitHub repository of the competition.

Benchmarks are not eligible for a prize, meaning that the total amount will be distributed among the competitors even if the benchmarks perform better than the forecasts submitted by the participants.

# M5

Similarly, any participating method associated with the organizers and the data provider, will not be eligible for a price.

## Submission

The forecasts for both competitions will be submitted through the Kaggle platform. The templates provided by the organizers through the platform can be used for this purpose.

두 경쟁에 대한 예측은 Kaggle 플랫폼을 통해 제출됩니다. 플랫폼을 통해 주최자가 제공 한 템플릿을 이 용도로 사용할 수 있습니다.

Note that the template of the point forecasts (M5 Forecasting - Accuracy) refers only to the 30,490 series that consist the lowest hierarchical level of the dataset (level 12 of Table 1) and not all 42,840 of the competition (all levels of Table 1). This is done because M5, in contrast to M4, M3, and other forecasting competitions where time series are mostly unrelated, deals among others with a real-life hierarchical forecasting problem. This means that the submitted forecasts must follow this hierarchical concept and, as a result, be coherent (forecasts at the lower levels have to sum up to the ones of the higher levels). In other words, it is assumed that the forecasting approach used for forecasting all 42,840 series of the competition derived coherent forecasts and, therefore, the forecasts of all levels can be automatically computed by aggregating (summing) the ones of the lowest level of the hierarchy.

포인트 예측 템플릿 (M5 예측-정확도)은 데이터 세트의 가장 낮은 계층 레벨 (표 1 의 레벨 12)을 구성하는 30,490 시리즈 만 참조하며 경쟁의 42,840 (표 1 의 모든 레벨)은 아닙니다. . M5 는 M4, M3 및 시계열이 거의 관련이없는 기타 예측 경쟁과 달리 실제 계층 적 예측 문제를 다루기 때문에 수행됩니다. 이는 제출 된 예측이이 계층 적 개념을 따라야하며 결과적으로 일관되어야합니다 (낮은 수준의 예측은 높은 수준의 예측과 합쳐 져야 함). 다시 말해, 42,840 개의 일련의 경쟁 파생 코 히어 런트 예측을 예측하는 데 사용되는 예측 방법이 사용되므로 모든 계층의 예측은 계층 구조의 가장 낮은 레벨의 집계를 합산하여 자동으로 계산할 수 있습니다..

It is important to note that the participants are completely free to use the forecasting approaches of their choice for forecasting the individual series. However, having done that, and by submitting just the forecasts of the lowest level, it will be assumed that the derived forecasts were reconciled before submitted for the final evaluation. For instance, a participant may forecast just the series at the bottom-level and derive the remaining forecasts using the bottom-up reconciliation method. Another participant may forecast just the series at the top level and get the ones at the lower levels using proportions (top-down reconciliation method). A mix of the previous two approaches is also possible (middle-out reconciliation method). Finally, predicting the series of all levels and getting the ones of the lowest level through an appropriate weighting scheme is also an option. The benchmarks describe some of these options, involving some indicative forecasting approaches that utilize the bottom-up (e.g. benchmark #12) and the top-down (e.g. benchmark #11) reconciliation method, as well as the combination of these two (e.g. benchmark #23).

# M5

참가자는 개별 시리즈를 예측하기 위해 자신이 선택한 예측 방법을 완전히 자유롭게 사용할 수 있습니다. 그러나 이를 수행하고 가장 낮은 수준의 예측만 제출하면 파생된 예측이 최종 평가를 위해 제출되기 전에 조정된 것으로 간주됩니다. 예를 들어, 참가자는 최저 수준의 계열만 예측하고 상향식 조정 방법을 사용하여 나머지 예측을 도출할 수 있습니다. 다른 참가자는 최고 수준에서 시리즈만 예측하고 비율을 사용하여 낮은 수준의 시리즈를 얻을 수 있습니다 (하향식 조정 방법). 이전 두 가지 접근 방식의 혼합도 가능합니다 (중간 조정 방법). 마지막으로, 적절한 가중치 체계를 통해 모든 수준의 시리즈를 예측하고 가장 낮은 수준의 수준을 얻는 것도 선택 사항입니다. 벤치마크는 상향식 (예: 벤치마크 # 12) 및 하향식 (예: 벤치마크 # 11) 조정 방법을 사용하는 일부 예측 예측 방법과 이 두 가지의 조합 (예: 벤치마크)을 포함하는 이러한 옵션 중 일부를 설명합니다. # 23).

Finally, given that there is not a direct and well-established way for reconciling probabilistic forecasts, the template of the probabilistic forecasts (M5 Forecasting - Uncertainty) requires inputting all 42,840 series of the competition. Thus, in this case, participants do not need to reconcile the forecasts using any of the above-mentioned approaches.

마지막으로, 확률 론적 예측을 조정하기 위한 직접적이고 잘 확립된 방법이 없기 때문에, 확률 론적 예측 (M5 예측-불확실성)의 템플릿은 42,840 시리즈의 경쟁을 모두 입력해야 합니다. 따라서 이 경우 참가자는 위에서 언급한 방법 중 하나를 사용하여 예측을 조정할 필요가 없습니다.