Evaluating generated in situ weather forecasts with exchangeability methods

David Landry, Anastase Charantonis, Claire Monteleoni INRIA Paris

Ver-AI workshop on Verification Reading 2025-06-23

Overview

Case study where multi-variate rank histogram methods helped in the evaluation and calibration of generative weather forecasting model.

- 1. Introduce an in situ generative postprocessing model
- 2. Diagnose a forecast spread issue using multivariate rank histograms
- 3. Calibrate the generative model

In situ forecasting, but the methodology should apply to gridded forecasts as well.

Generative in situ forecasting model

Flow matching postprocessing (FMAP)

Generative in situ postprocessing method with a particular focus modeling on **spatial correlations** (Landry et al. 2025).

Based on **flow matching** (very closely related to **denoising diffusion methods**).



Key characteristics

The statistical objects the generative model approximates are **samples** (as opposed to the conditional expectation).

A numerical integration displaces the samples from a normal distirbution to the target distribution.

$$p_0(x)$$

$p_0({\boldsymbol{x}}) = \mathcal{N}({\boldsymbol{0}}, {\boldsymbol{I}})$

 $p_1(x) \approx \text{target distribution}$

Flow matching distribution transport



Example flow integration

Experiments

Dataset EUPPBench (Demaeyer et al. 2023)

- 0.25° resolution
- 20 years reforecasts (bi-weekly)
- 2 years forecasts (daily)

Target: Surface temperature and wind gust at 121 stations in central Europe.



Flow matching model skill score



Spread-error ratio for flow matching model (FMAP) and other baselines

Dispersion behavior of generative models

Both **GenCast** (Price et al. 2025) and **ArchesWeatherGen** (Couairon et al. 2024) are under-dispersive.

The latter proposes a calibration procedure to compensate.

Multivariate rank histograms to evaluate dispersion

Evaluation strategy

We calibrate the well-known rank histograms, but in a multivariate setting.

Our multivariate vector x is the forecast for one variable at every station (121-long).

Evaluation strategy

We will use methods based on

- Minimum Spanning Tree (MST) (Smith and Hansen 2004; Wilks 2004)
- \cdot The Mahalanobis distance

These methods give us some insight on the dispersion behavior of the models.

Exchangeability methods for evaluation

We assume the ensemble

 $\{m{x}_0,m{x}_1,...,m{x}_m\}$

is **exchangeable** where x_0 is the verifying observation and $x_{1...m}$ the ensemble members.

Informally, we verify that observation x_0 behaves like any other member.

Pre-rank functions

A straightforward strategy to build rank histograms for multivariate quantities is to use a **pre-rank function** (Gneiting et al. 2008).

$$F: \mathbb{R}^d \times \underbrace{\mathbb{R}^d \times \ldots \times \mathbb{R}^d}_{\text{Symmetric}} \longrightarrow \mathbb{R}$$

The prerank of member j is

$$z_j = F(\boldsymbol{x}_j; \boldsymbol{x}_{-j})$$

where $x_{-j} = \{x_{0..m}\} \setminus ig\{x_jig\}.$

Allowable pre-rank functions are **symmetric** on the right-hand side argument to preserve exchangeability (rank histogram flatness).

Minimum spanning tree

A proposal by (Smith and Hansen 2004; Wilks 2004)

Build a fully-connected graph where

- the nodes are the ensemble members
- the edges have length $\| \boldsymbol{x}_i \boldsymbol{x}_j \|$

From that graph extract the **minimum spanning tree** (MST) and measure its length.

$$F_{\mathrm{MST}(\boldsymbol{x}_{j};\boldsymbol{x}_{-j})} = \mathrm{MST}(\boldsymbol{x}_{-j})$$

Interpreting MST histograms



When the observation is removed, the MST is shorter. Underdispersion, systematic bias.

Interpreting MST histograms



When the observation is removed, the MST is longer. Overdispersion. $_{1}$

Mahalanobis distance

Sample to normal distribution metric.

Given a distribution $Q = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ we have

$$\mathrm{MD}(Q, \boldsymbol{y}) = \sqrt{(\boldsymbol{\mu} - \boldsymbol{y})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu} - \boldsymbol{y})}.$$

Mahalanobis prerank function

We propose

$$z_j = F_{\mathrm{MD}}(\boldsymbol{x}_j; \boldsymbol{x}_{-j}) = \mathrm{MD}(Q_{\boldsymbol{x}_{-j}}, \boldsymbol{x}_j)$$

where $Q_{x_{-j}}$ is estimated from the ensemble. This may require specific estimation methods such as a Ledoit-Wolf shrinkage.

Has a similar interpretation to the MST (in reverse).

Interpreting MD histograms



Observation is far from the ensemble distribution \rightarrow Underdispersion, systematic biases. 22

Interpreting MD histograms



Observation is central within the ensemble \rightarrow Overdispersion.

What about the Box Ordinate Transform (BOT)?

This proposed method is similar to the BOT evaluation method (Gneiting et al. 2008). Instead of a prerank it computes

$$u = 1 - \chi_d^2 \left[\left(\boldsymbol{x}_0 - \boldsymbol{\mu} \right)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x}_0 - \boldsymbol{\mu}) \right]$$

which we translate to a bin number in a PIT histogram using

$$b = 1 + \lfloor u(m+1) \rfloor.$$

We get uniform u if the normality assumption is true.

What about the Box Ordinate Transform (BOT)?

Our model fails the BOT test consistently

• Normal assumption is false – the χ^2 test saturates

Our Mahalanobis pre-rank approach is less sensitive to the normality of the data. It ensures that all members are equally non-normal.

It requires m times more matrix inversions.

Experiments on EUPP postprocessing

Exchangeability test. Prerank: Minimum Spanning Tree. Lead time: 3D.



MST test on 84 days calibration period

Exchangeability test. Prerank: Mahalanobis. Lead time: 3D.



Mahalanobis distance test on 84 days calibration period

Calibrating the flow matching model

Calibrating the flow matching model

We propose to use the multivariate rank histograms to calibrate the flow matching model as in Couairon et al. (2024).

The initial sample is rescaled such that

 $p_0(\pmb{x}) = \mathcal{N}(\pmb{0}, \alpha \pmb{I}).$

We reserved an 84 days calibration period (1st week of each month in 2017).

Exchangeability test. Prerank: Minimum Spanning Tree. Lead time: 3D.



Rank histograms for flow matching model calibration (MST)

Exchangeability test. Prerank: Mahalanobis. Lead time: 3D.



Rank histograms for flow matching model calibration (Mahalanobis)

Flow matching calibration Lead time: 3D



Chi-square test statistic given the scaling factor

Flow matching calibration Lead time: 3D



Chi-square test statistic given the scaling factor

Flow matching calibration



Stacked all lead-time-variable combinations



This calibration would require a larger calibration set using the ES

Effect of calibration on other metrics

Flow matching Energy Skill Score after calibration







Discussion

The α parameter has an important impact on dispersion, as measured by out reliability tests.

Calibration seems successful for temperature, but less so on wind gust.

Any interactions with the normalization procedure on wind?

Conclusion

In summary...

- Multivariate rank histograms were useful in evaluating an in situ generative weather forecasting model
- We proposed an exchangeability test similar to BOT which is less sensitive to the normality of the data
- This allowed a multivariate calibration procedure would have been challenging on the Energy Score

Conclusion

In the future...

- How does this translate to higher dimensionalities like in full NWP forecasts?
- Better account for test multiplicity and serial correlations (Wilks 2019)

References

- Couairon, G., R. Singh, A. Charantonis, C. Lessig, and C. Monteleoni, 2024: ArchesWeather & ArchesWeatherGen: A Deterministic and Generative Model for Efficient ML Weather Forecasting. https://doi.org/10.48550/arXiv.2412.12971.
- Demaeyer, J., and Coauthors, 2023: The EUPPBench Postprocessing Benchmark Dataset v1.0. *Earth System Science Data Discussions*, 1–25, https://doi.org/10. 5194/essd-2022-465.
- Gneiting, T., L. I. Stanberry, E. P. Grimit, L. Held, and N. A. Johnson, 2008: Assessing Probabilistic Forecasts of Multivariate Quantities, with an Application to Ensemble Predictions of Surface Winds. *TEST*, **17**, 211–235, https://doi.org/10.1007/s11749-008-0114-x.
- Landry, D., C. Monteleoni, and A. Charantonis, 2025: Generating Ensembles of Spatially-Coherent in-Situ Forecasts Using Flow Matching. https://doi.org/10. 48550/arXiv.2504.03463.
- Price, I., and Coauthors, 2025: Probabilistic Weather Forecasting with Machine Learning. Nature, 637, 84–90, https://doi.org/10.1038/s41586-024-08252-9.
- Smith, L. A., and J. A. Hansen, 2004: Extending the Limits of Ensemble Forecast Verification with the Minimum Spanning Tree. *Monthly Weather Review*, **132**, 1522–1528.
- Wilks, D. S., 2004: The Minimum Spanning Tree Histogram as a Verification Tool for Multidimensional Ensemble Forecasts. *Monthly Weather Review*, **132**, 1329–1340.
- Wilks, D. S., 2019: Statistical Methods in the Atmospheric Sciences: An Introduction. 4th ed. Elsevier,.

Thank you





FMAP preprint

Email