

SPLICED BINNED-PARETO DISTRIBUTION FOR ROBUST MODELING OF HEAVY-TAILED TIME SERIES

Elena Ehrlich
 AWS ProServe
 Miami, FL, USA
 eeehrlic@amazon.com

Laurent Callot
 Amazon Research
 Seattle, WA, USA
 lcallot@amazon.com

François-Xavier Aubet
 Amazon Research
 Vienna, Austria
 aubetf@amazon.com

ABSTRACT

This work proposes a novel method to robustly and accurately model time series with heavy-tailed noise, in non-stationary scenarios. In many practical application time series have heavy-tailed noise that significantly impacts the performance of classical forecasting models; in particular, accurately modeling a distribution over extreme events is crucial to performing accurate time series anomaly detection. We propose a Spliced Binned-Pareto distribution which is both robust to extreme observations and allows accurate modeling of the full distribution. Our method allows the capture of time dependencies in the higher order moments of the distribution such as the tail heaviness. We compare the robustness and the accuracy of the tail estimation of our method to other state of the art methods on Twitter mentions count time series.

1 INTRODUCTION

In many real world applications, time series can have heavy-tailed distributions, examples include: (i) financial series where speculators who bet high, bet very high (Bradley & Taquq, 2003), (ii) server metrics like requests-per-second to plan compute resource scaling, (iii) extreme rainfall for flood prediction (Bezak et al., 2014). Extreme events are known to impede estimation of the predictive density $p(\mathbf{x}_{t+1}|\mathbf{x}_{1:t})$ for a time series $\mathbf{x}_{1:T}$, $\mathbf{x}_t \in \mathcal{X} \subset \mathbb{R}$, $t, T \in \mathbb{N}_1$, under classical methods. This problem has been extensively studied in the case of independently and identically distributed (iid) data (Beirlant et al., 2006), and has recently gained attention in the context of time series (Kulik & Soulier, 2020; Davison & Huser, 2015). Accurate estimation of distribution tails can be important in forecasting settings and is crucial in time series anomaly detection, where how likely the extreme events are dictates whether an alarm should be raised.

Robust and accurate density estimation is challenging in time series with heavy-tailed noise for two main reasons: first, extreme events have a strong impact on the estimation of the base distribution; second, one has to account for time-varying components of the time series in order to identify the tails before beginning to fit the tails themselves. While some methods have been proposed for tail estimation, they assume little to no time-varying components Siffer et al. (2017). To obtain a forecast robust to extreme values and adjustable to many shapes of distributions, discrete binned distributions can be used Rabanser et al. (2020), these are parametrised by a Neural Network (NN), which have been shown capable of capturing complex long-term time dependencies in forecasting (Benidis et al., 2020), We address these two challenges by combining a binned distribution with Generalised Pareto distributions for the tails, all three distributions parameterised by a single NN, allowing us to jointly model time dependencies in the base distribution and in the tails.

In particular, we make the following contributions:

1. We propose a combination of a binned distribution, for a robust estimate of the base distribution, with Generalised Pareto distributions, for an accurate estimate of the tails.
2. Our method allows for asymmetric tails with time-varying heaviness and scale.
3. We show empirically on a real world example that our approach allows for a better estimate of the tail distribution than previous methods, while also providing an accurate estimate of the base distribution.

2 BACKGROUND & RELATED WORK

Extreme Value Theory (EVT) A classic result from extreme value theory (EVT) states that the distribution of extreme values is almost independent of the base distribution of the data (Fisher & Tippett, 1928). As a consequence, it has been proposed to estimate only the tail of the distribution by considering only the peaks above threshold $\tau \in \mathcal{X}$, the assumed upper bound of the base distribution. Let $\bar{F}_\tau(x) = P(X > x + \tau | X > \tau)$; by the second theorem of EVT (Balkema & De Haan, 1974; Pickands III et al., 1975):

$$\bar{F}_\tau(x) \sim \text{GPD}(\xi, \beta) \quad (1)$$

where GPD is a Generalised Pareto Distribution with shape $\xi \in \mathbb{R}$ and scale $\beta \in \mathbb{R}^+$. Therefore the quantile value $z_q : P(X > z_q) < q$ for quantile level $q \in [0, 1]$ can be calculated

$$z_q \approx \tau + \frac{\hat{\beta}}{\hat{\xi}} \left(\left(\frac{qT}{N_\tau} \right)^{-\hat{\xi}} - 1 \right) \quad (2)$$

upon solving for $\hat{\xi}, \hat{\beta}$ and where N_τ is the number of peaks over initial threshold τ out of T observations. This method of computing the quantile of extreme values is known as Peaks-Over-Threshold (POT) (Leadbetter, 1991). Note that Eq 1 only holds as τ tends to the true upper bound of the base distribution. In practice τ is not known, and has to be determined. This introduces a trade-off between the variance incurred from too few excesses above an over-estimated threshold and the bias incurred from non-tail observations introduced into the GPD fit.

Streaming POT and Drift SPOT Siffer et al. (2017) extended the result of POT to the streaming setting and for incorporation of online anomaly detection in stationary and concept-drift time series, known respectively as Streaming POT (SPOT) and Drift SPOT (DSPOT). They aim to accurately estimate a high quantile z_q above which they consider incoming points as anomalous. Using a fixed threshold τ , they use POT to initialise z_q from the first T observations. The method then processes data arrivals of $X_t < \tau$ as ‘normal’, $X_t > z_q$ as anomalous, and $\tau \leq X_t \leq z_q$ as non-anomalous ‘real’ peaks that trigger a recalculation of Eq 1-2 to update z_q .

Slightly relaxing the stationarity restriction of SPOT, DSPOT subtracts a sliding window average of ‘normal’ points before applying SPOT. Though adaptable to slow drifts in the distribution mean, DSPOT remains limited in three ways: 1) it does not model complex time-dependencies in the base distribution, 2) it does not model the base distribution and is therefore not suitable for forecasting, 3) it does not allow for a time-varying, feature dependent parameterisation of the tails.

NN extensions of SPOT To further improve DSPOT’s relaxation of the stationarity constraint Davis et al. (2019) proposed point forecasting by training an RNN to minimise the Mean Squared Error (MSE), and then using SPOT on the prediction residuals to detect anomalies. The solution is limited in three ways: 1) the MSE loss function is sensitive to extreme events, 2) no allowance is made for time dependencies in moments of the distribution beyond the mean, be it the variance or the tail heaviness, 3) the use of residuals prohibits modelling asymmetric tails.

3 SPLICED BINNED-PARETO DISTRIBUTION

We propose the Spliced Binned-Pareto (SBP) distribution which uses a flexible binned distribution to model the base of the distribution and two GPDs to model the tails.

We model the base of the predictive distribution with a discrete binned distribution to make it robust to extreme values and adaptable to the variety of real work distributions. As described by Rabanser et al. (2020), we discretise the real axis between two points into n bins. A NN is trained to predict the probability of the next point falling in each of these bins, as shown in Figure 1a. This gives a distribution robust to extreme values at training time because it is now a classification problem, the log-likelihood is not affected by the distance between the predicted mean the observed point, as would be the case when using a Gaussian or Student’s-t distribution for example.

We enhance the binned distribution to have an accurate parametric estimate of the tails by replacing the tails of the binned with GPDs. From the binned distribution, we delimitate the base distribution

from its upper and lower tails according to user-defined quantiles q and $1 - q$ respectively (e.g. $q = 0.05$). At time t we obtain $\tau_t^{\text{lower}} = z_{1-q}(t)$ from the predictive binned distribution and replace the cumulative distribution function (cdf) of the bins below this quantile with the weighted cdf of the lower GPD. Analogously we obtain $\tau_t^{\text{upper}} = z_q(t)$ and replace the cdf above with the weighted upper GPD. Using this procedure we obtain a valid probability distribution, integrating to 1, with support $(-\infty, \infty)$ (illustration in Figure 1b).

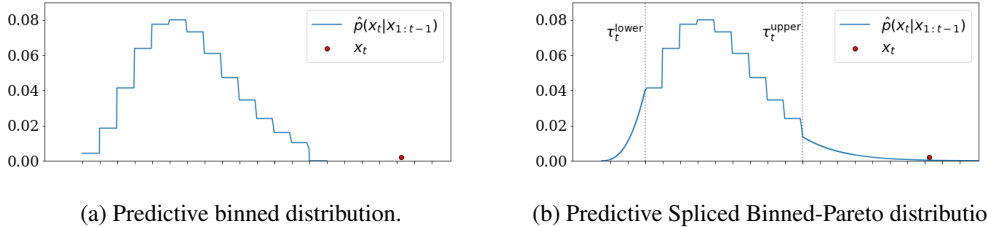


Figure 1: Illustration of the parametric tail accounting for an extreme event.

The binned distribution and the GPDs are parameterised by a single NN taking as input $\mathbf{x}_{1:t-1}$, the past points, and outputting $n + 4$ parameters: each of the n bin probabilities as well as ξ and β for each of the GPD. While our approach is fully general and works with most NN architectures, we opt to use a Temporal Convolution Network (TCN) (Bai et al., 2018). We expect our results to carry over to other architectures.

In addition to being robust to extreme values, it also results in a robust estimate of the tails of the distributions as it can model changes in the tails over time. While predecessor methods use a fixed threshold τ to delimit the tails, by modeling the base distribution we obtain a time-varying threshold. Furthermore, training a single NN parameterising all distributions to maximise the log probability of the observed time step under the binned and GPD distributions, results in an prediction that accounts for temporal variation in all moments of the distribution: the mean and variance as well as tail heaviness and scale. This includes asymmetric tails.

4 EXPERIMENTS

We compare our approach to the different methods presented in Related Work: SPOT, DSPOT, as well as SPOT on NN prediction residuals. We implemented Davis et al. (2019) using a TCN to restrict ourselves to differences in the extreme value handling; we refer to it as TCN-SPOT. Following (Siffer et al., 2017), for each of SPOT, DSPOT and TCN-SPOT, we set threshold τ at $z_{0.95}$ of the training data for the lower tail and at $z_{0.05}$ for the upper tail. And we set the thresholds of the SBP with $q = 0.05$, replacing the lower and upper 5% of the distribution with GPDs.

Evaluation metric We evaluate the accuracy of the density estimation of each of the method using Probability-Probability plots (PP-plots) Michael (1983). For a given quantile level q , we compute y_q the fraction of points that fell below the given quantile $z_q(t)$ of their corresponding predictive distribution:

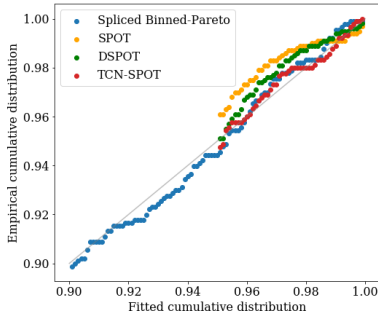
$$y_q = \frac{\sum_{t=2}^T \mathbb{I}[\mathbf{x}_t < z_{1-q}(t)]}{T}, \quad z_q(t) : p(\mathbf{x}_t > z_q(t) | \mathbf{x}_{1:t-1}) < q \quad (3)$$

To obtain a quantitative score, we measure how good the tail estimate is by computing the Mean Absolute Error (MAE) between y_q and q for all measured quantiles q .

4.1 EVALUATION ON SYNTHETIC DATA

To compare the different methods on simple time dependencies, we generated sine waves and added iid Student’s-t and heavy-tailed noises for synthetic data. Figure 2a shows the PP-plot for each comparative method. We observe that our method provides an accurate estimate of the density both of the top of the base distribution, $\text{cdf} \in [0.90, 0.95]$, and of the tail of the distribution, $\text{cdf} \in [0.95, 1.0)$, and does so without discontinuity at $\tau_t = z_{0.05}(t)$. Further, while one could have feared

that parametrising the GPD using a NN could have been noisy or less robust, in fact our method obtains a better estimate of the tail than SPOT or DSPOT.



(a) PP-plot on the sinwave time series for $cdf \in [0.9, 1.0)$.

Model	Synthetic data	Real data
SPOT	9.0×10^{-3}	9.61×10^{-3}
DSPOT	5.65×10^{-3}	1.08×10^{-2}
TCN-SPOT	2.06×10^{-3}	8.77×10^{-3}
SBP (ours)	3.17×10^{-3}	4.24×10^{-3}

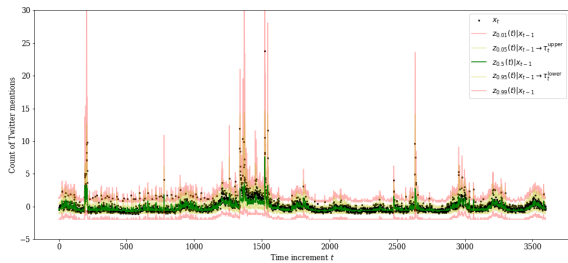
(b) Mean absolute error between the empirical quantiles and the predicted quantiles, on the synthetic and real datasets. (lower error is better)

Figure 2: Evaluation on synthetic time series and comparison of the MAE .

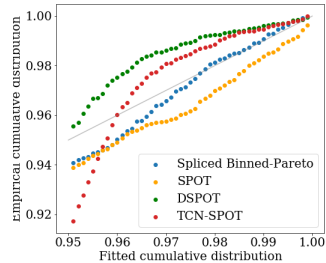
4.2 EVALUATION ON A REAL WORLD SCENARIO

Retailers rely upon demand forecasting to inform inventory ordering. Recently, the reach of social media platforms has meant that social trends take off faster than the forecast can anticipate. While the first extreme demand spike is not predictable, it is still important to reliably quantify its likelihood. To illustrate this point, we look at the count of Twitter mentions that a stock ticker symbol of interest receives per 5-minute interval (data source: Numenta Ahmad et al. (2017)). This series exhibits seasonalities, heteroscedasticity, and extreme realizations; see Figure 3a.

Figure 3b shows the PP-plots on the real time series, focusing on the tail estimate for $cdf \in [0.95, 1.0)$. We observe that our method obtains a significantly better tail estimate than the comparison partners. Table 2b shows that, while in the setting with simple time dependencies and symmetric tail heaviness our method performs comparably to other method, in the real scenario the advantage of our method’s robustness becomes clear.



(a) Spliced Binned-Pareto fit on the test set of real time series.



(b) PP-plot of the upper tail.

Figure 3: Time series of the stock ticker’s Twitter mentions per 5-minute interval

5 DISCUSSION

This work presents the Spliced Binned-Pareto Distribution distribution which, combined with a TCN, allows robust and accurate estimation of the predictive density in the presence of extreme events. The bias variance trade-off inherent to setting τ in SPOT is less present in our method as the thresholds are time-varying; however our method still requires a user-defined quantile to delimitate each tail. We want to investigate approaches to learn the tail quantile from training data to further reduce the domain knowledge needed to use the method.

REFERENCES

- Subutai Ahmad, Alexander Lavin, Scott Purdy, and Zuha Agha. Unsupervised real-time anomaly detection for streaming data. *Neurocomputing*, Available online 2 June 2017, ISSN 0925-2312, <https://doi.org/10.1016/j.neucom.2017.04.070>. 262:134–147, 2017.
- Shaojie Bai, J Zico Kolter, and Vladlen Koltun. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271*, 2018.
- August A Balkema and Laurens De Haan. Residual life time at great age. *The Annals of probability*, pp. 792–804, 1974.
- Jan Beirlant, Yuri Goegebeur, Johan Segers, and Jozef L Teugels. *Statistics of extremes: theory and applications*. John Wiley & Sons, 2006.
- Konstantinos Benidis, Syama Sundar Rangapuram, Valentin Flunkert, Bernie Wang, Danielle Maddix, Caner Turkmen, Jan Gasthaus, Michael Bohlke-Schneider, David Salinas, Lorenzo Stella, et al. Neural forecasting: Introduction and literature overview. *arXiv preprint arXiv:2004.10240*, 2020.
- Nejc Bezak, Mitja Brilly, and Mojca Šraj. Comparison between the peaks-over-threshold method and the annual maximum method for flood frequency analysis. *Hydrological Sciences Journal*, 59(5):959–977, 2014.
- Brendan O Bradley and Murad S Taqqu. Financial risk and heavy tails. In *Handbook of heavy tailed distributions in finance*, pp. 35–103. Elsevier, 2003.
- Neema Davis, Gaurav Raina, and Krishna Jagannathan. Lstm-based anomaly detection: Detection rules from extreme value theory. In *EPIA Conference on Artificial Intelligence*, pp. 572–583. Springer, 2019.
- Anthony C Davison and Raphaël Huser. Statistics of extremes. *Annual Review of Statistics and its Application*, 2:203–235, 2015.
- Ronald Aylmer Fisher and Leonard Henry Caleb Tippett. Limiting forms of the frequency distribution of the largest or smallest member of a sample. In *Mathematical Proceedings of the Cambridge Philosophical Society*, volume 24, pp. 180–190. Cambridge University Press, 1928.
- Rafal Kulik and Philippe Soulier. *Heavy-tailed time series*. Springer, 2020.
- M.R. Leadbetter. On a basis for ‘peaks over threshold’ modeling. *Statistics & Probability Letters*, 12(4):357 – 362, 1991. ISSN 0167-7152. doi: [https://doi.org/10.1016/0167-7152\(91\)90107-3](https://doi.org/10.1016/0167-7152(91)90107-3).
- John R Michael. The stabilized probability plot. *Biometrika*, 70(1):11–17, 1983.
- James Pickands III et al. Statistical inference using extreme order statistics. *the Annals of Statistics*, 3(1):119–131, 1975.
- Stephan Rabanser, Tim Januschowski, Valentin Flunkert, David Salinas, and Jan Gasthaus. The effectiveness of discretization in forecasting: An empirical study on neural time series models. *arXiv preprint arXiv:2005.10111*, 2020.
- Alban Siffer, Pierre-Alain Fouque, Alexandre Termier, and Christine Largouet. Anomaly detection in streams with extreme value theory. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1067–1075, 2017.