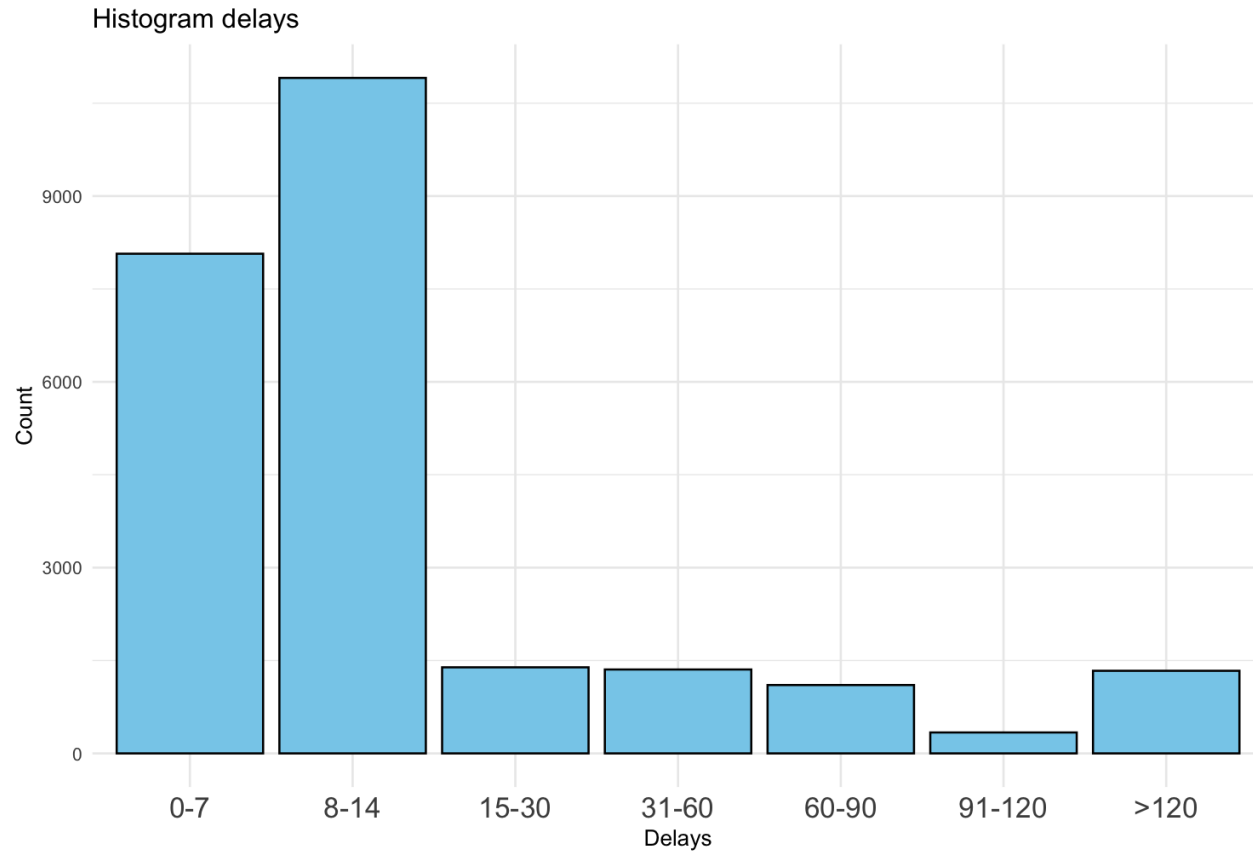# Presentation

# GLOBAL VIEW : DESCRIPTIVE
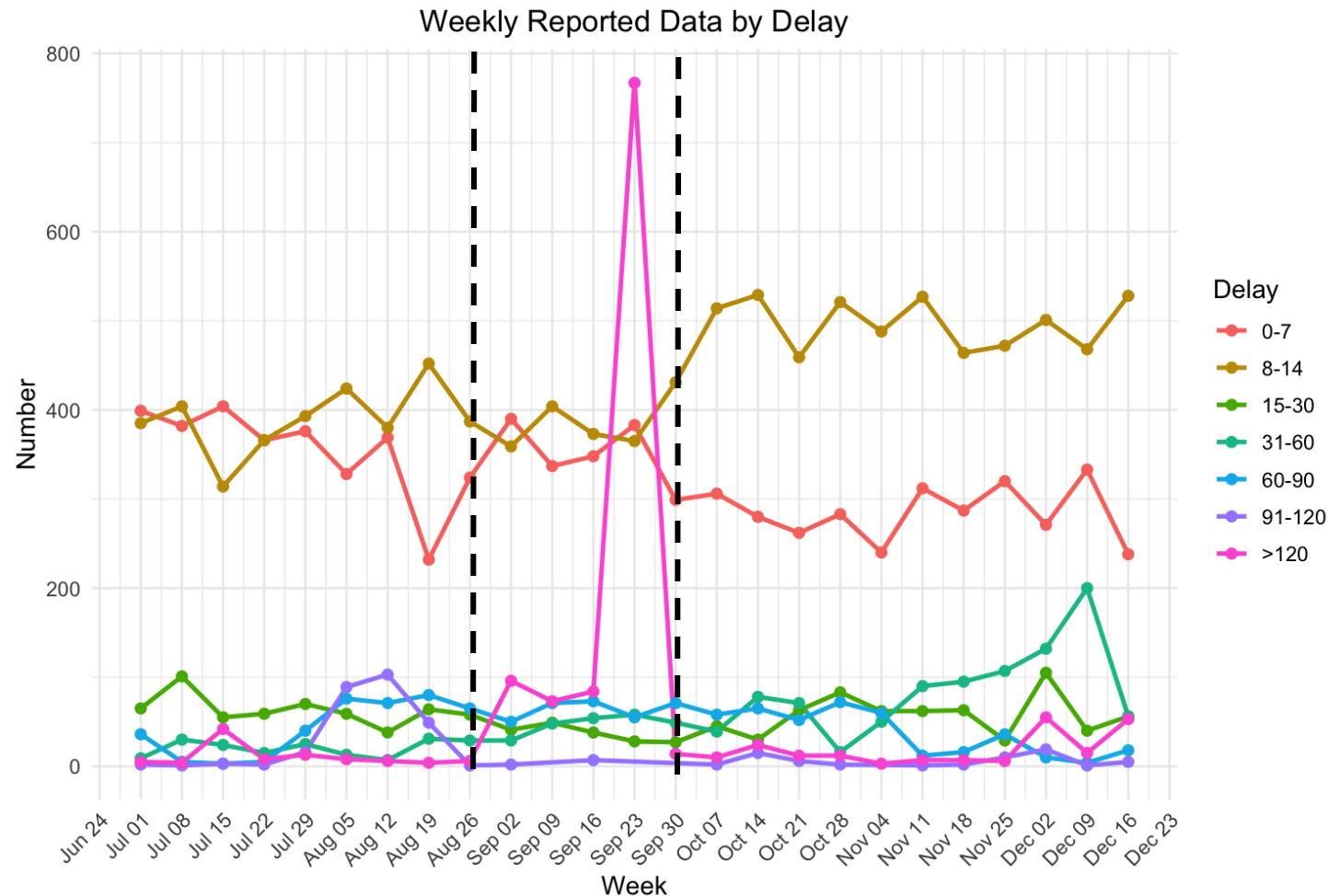


Histogram delays

Data: the events that newly appeared after the first download are assumed to be newly reported

**Delay: date of first report - date of event**

- The date of report is subject to change when there is an update (sometimes an error correction)

- For the global view of the delay we consider the first date of report and omit error correction (this comes later)

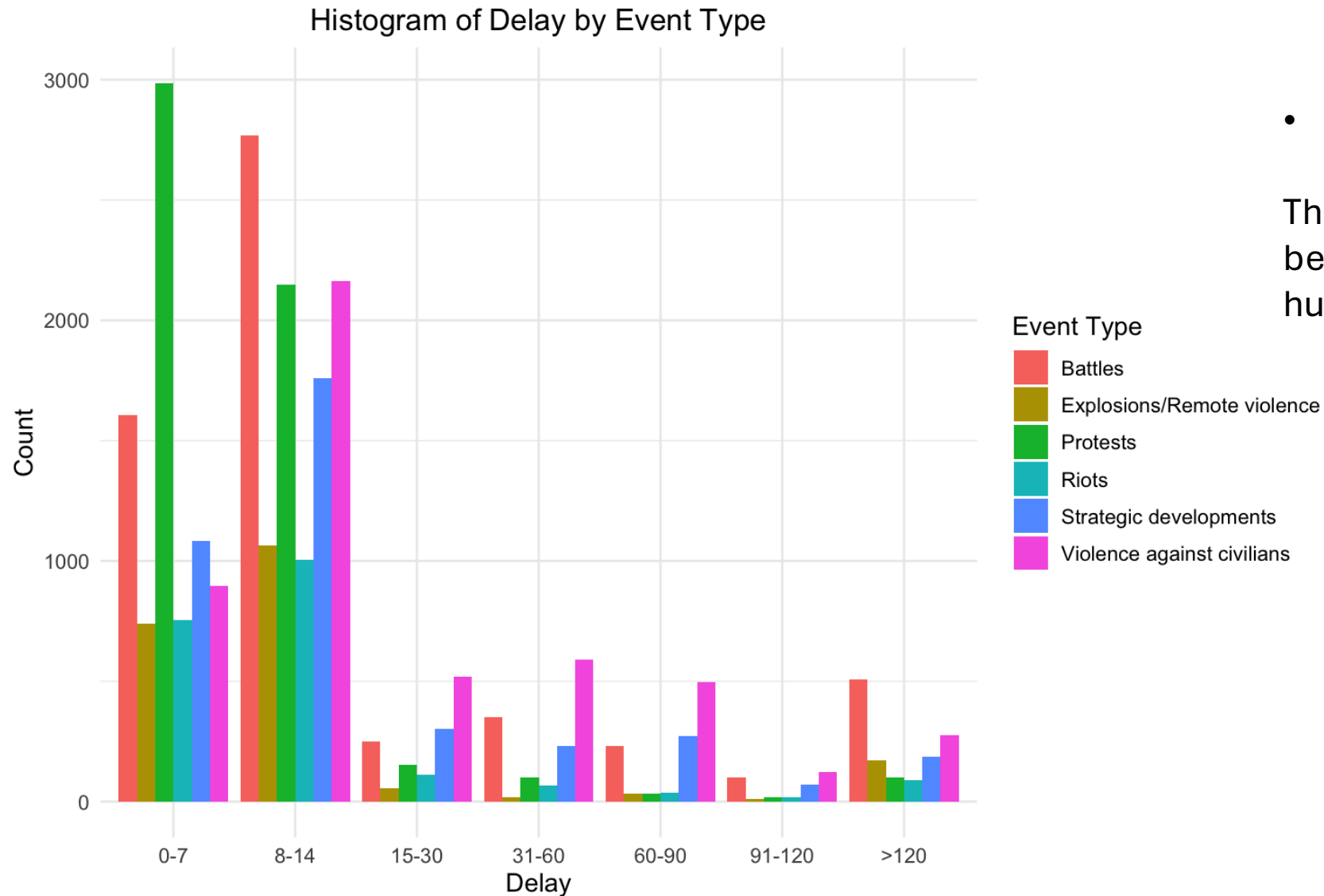# GLOBAL VIEW : DESCRIPTIVE



Weekly Reported Data by Delay

Number of events reported weekly by delay

| country | Count |
|:-------------|-----:|
| Cameroon | 3 |
| Somalia | 759 |
| South Africa | 3 |
| South Sudan | 2 |

We have sko symmetry between 0-7 and 8-14 (this indicates a negative correlation)

**Corr = -0.8**
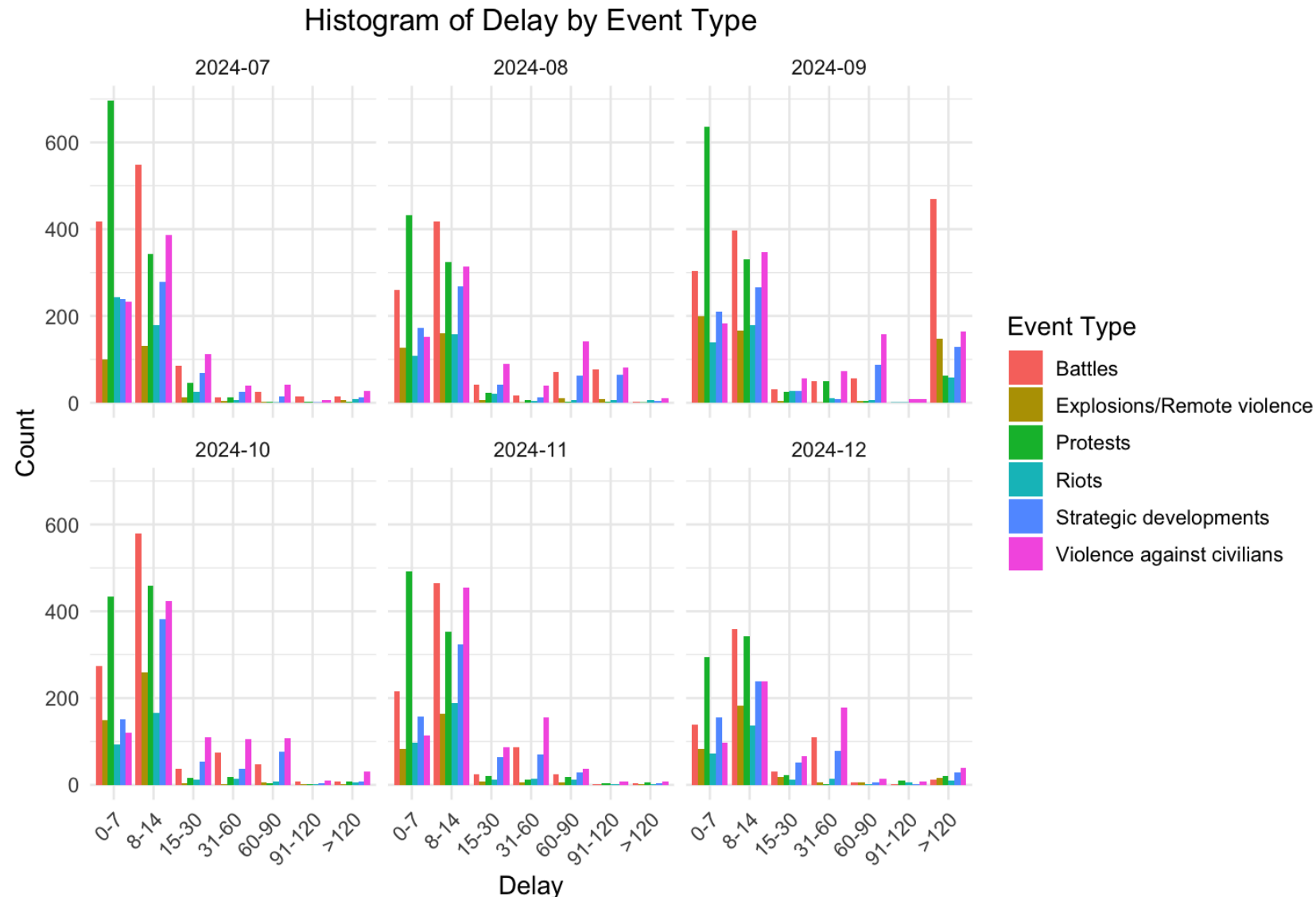
**Are the delays autocorrelated? (and of what order)**

# GLOBAL VIEW : DESCRIPTIVE



Histogram of Delay by Event Type

- How delays are distributed between event types?

- Does it change over time?

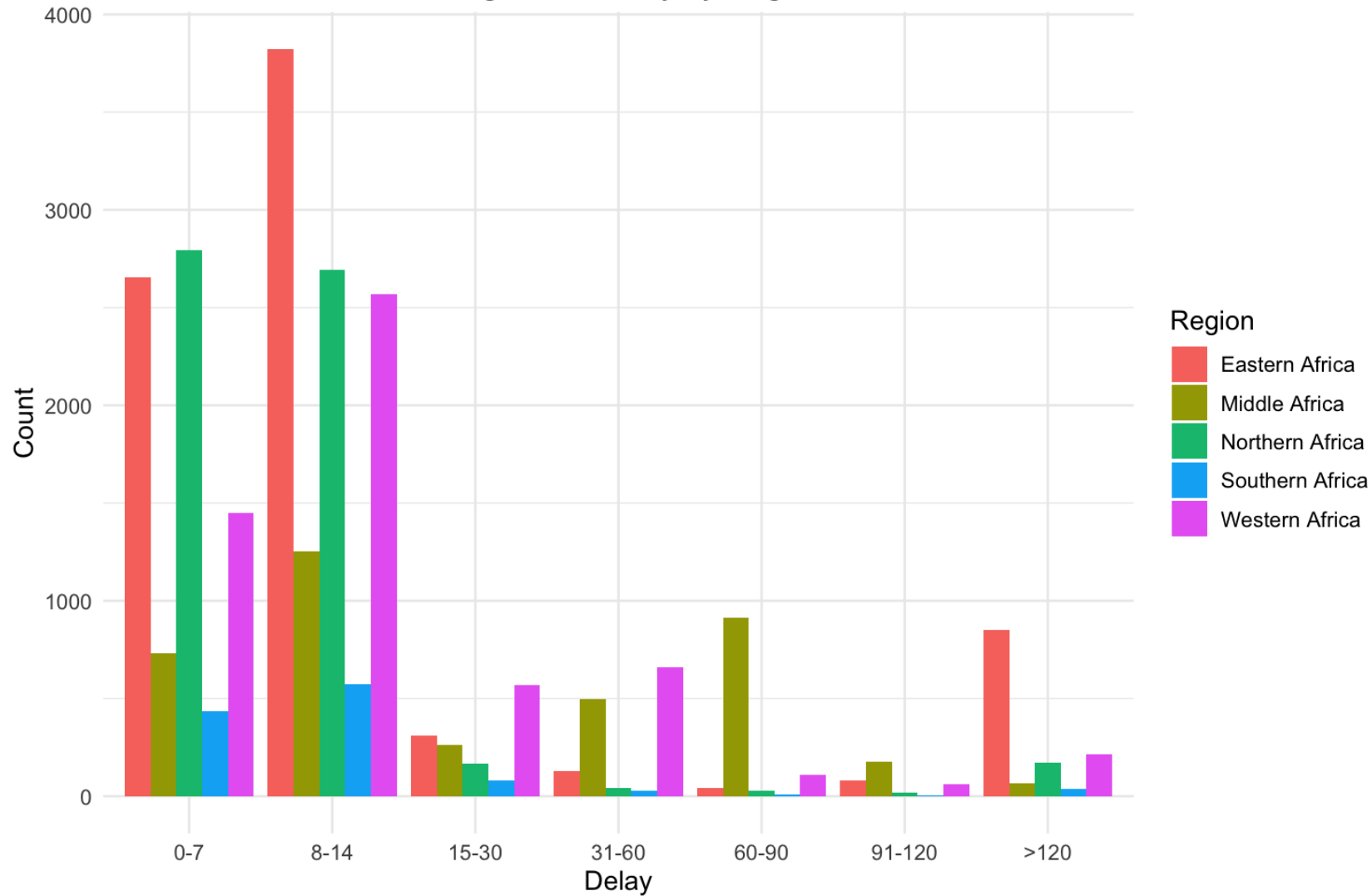These questions will be answered better with KM but some histo won't hurt

# GLOBAL VIEW : DESCRIPTIVE



Histogram of Delay by Event Type

- The effect of **event type** on delay does not vary much over time

- Can we assume a constant effect? **We'll find out more with the KM**

- **What happened in sept 2024. Very different from the other months in the >120**
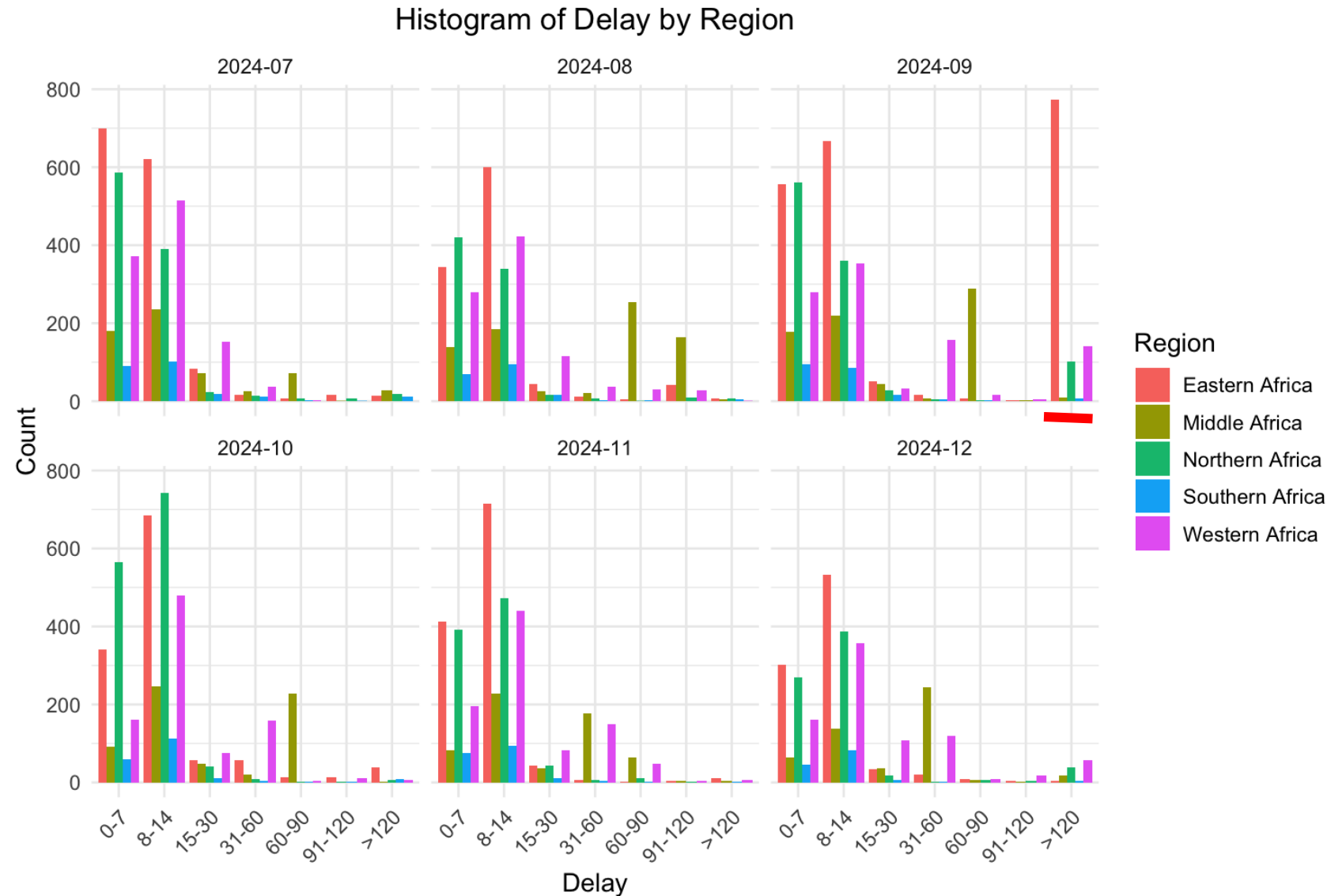
# GLOBAL VIEW : DESCRIPTIVE
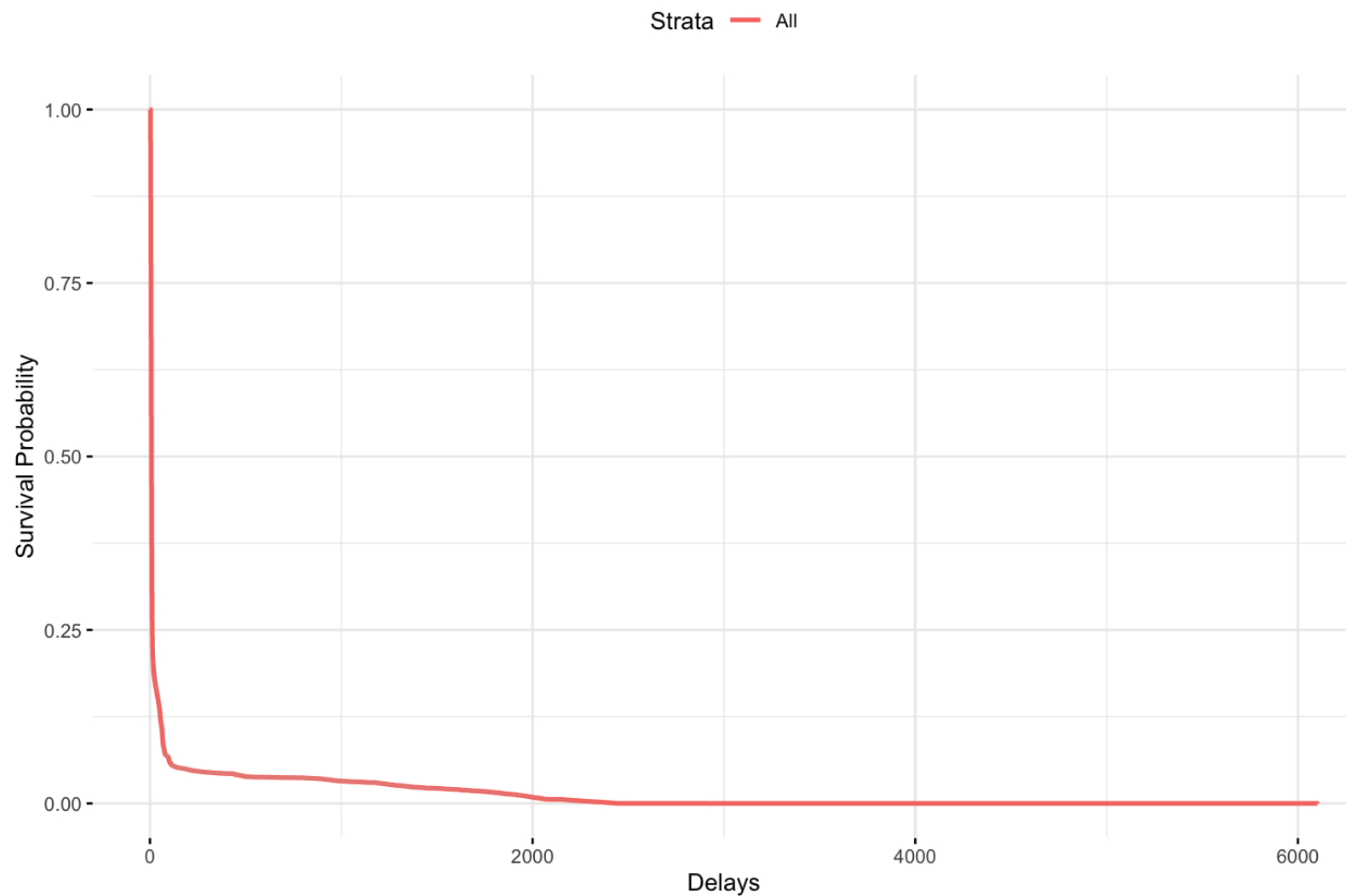

Histogram of Delay by Region

- How delays are distributed between regions?

- **Regions where there are many conflicts get report more on time?**

# GLOBAL VIEW : DESCRIPTIVE


Histogram of Delay by Region

- The effect of **region** on delay does not vary much over time

- Can we assume a constant effect? **We'll find out more with the KM**
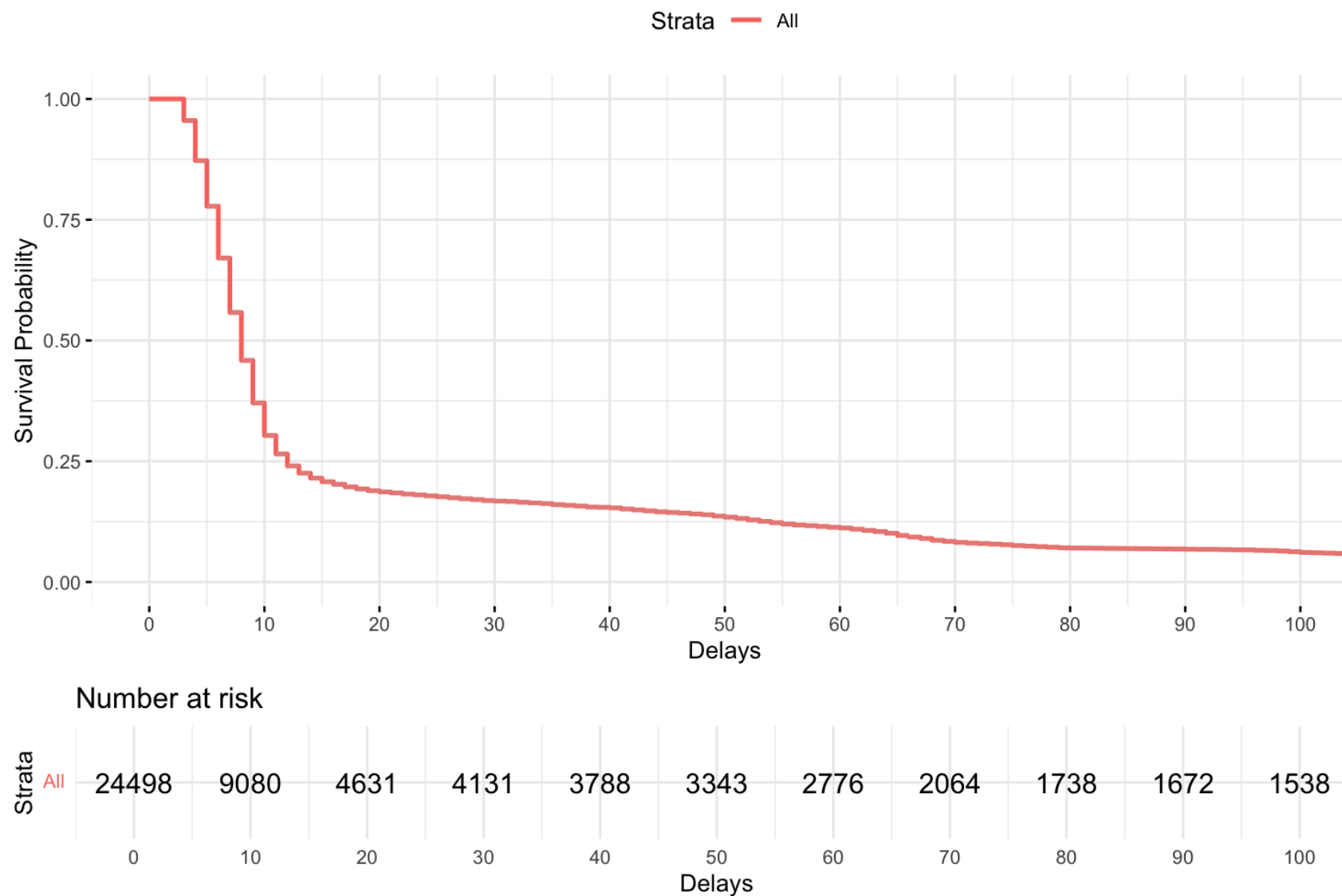
# KM



Long tail but where do we right truncate it?

The survival function decreases considerably between 0 and 80.
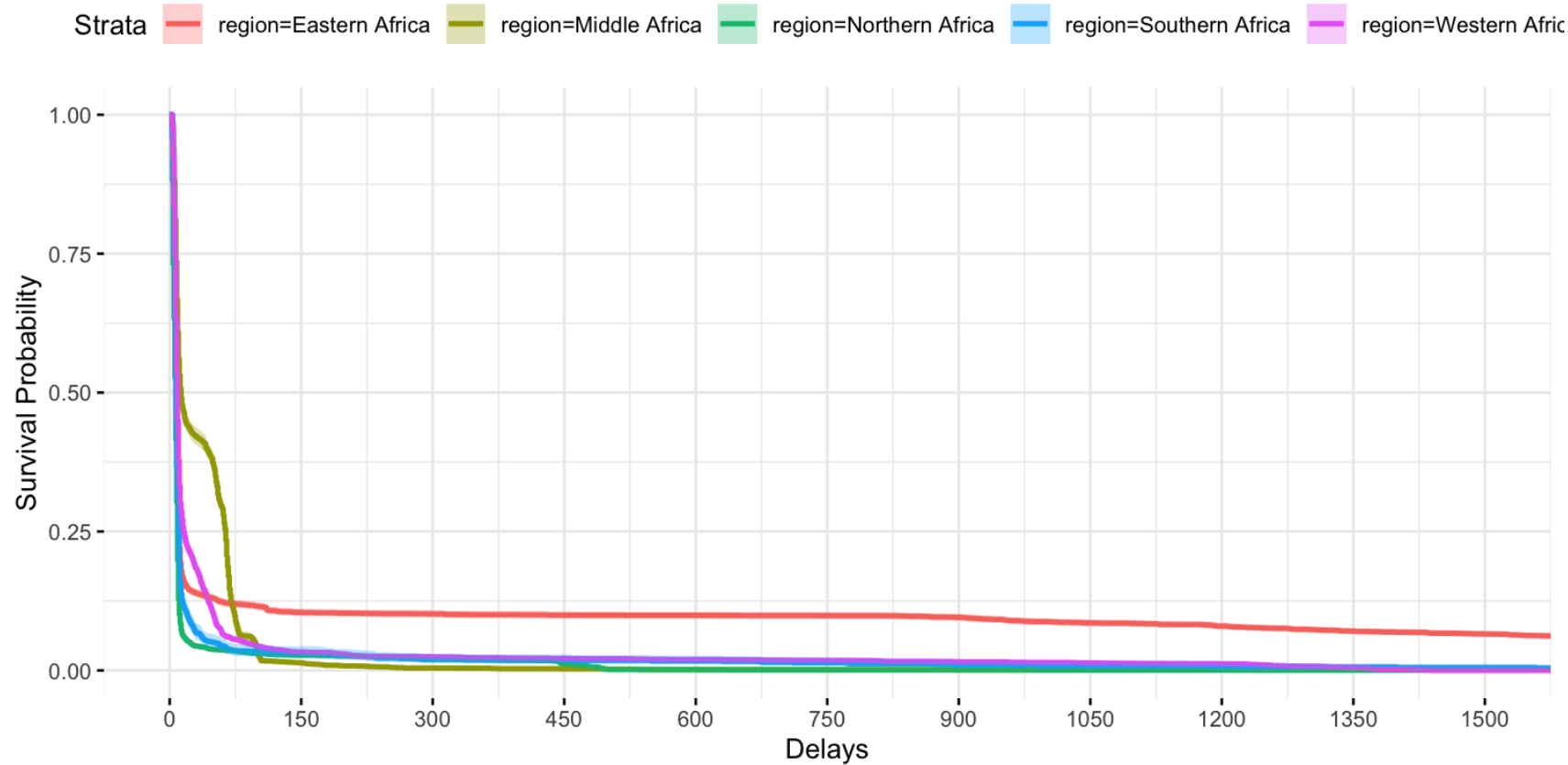
Short term delays are relatively more predominant

# KM



S_hat(100) = 0.063

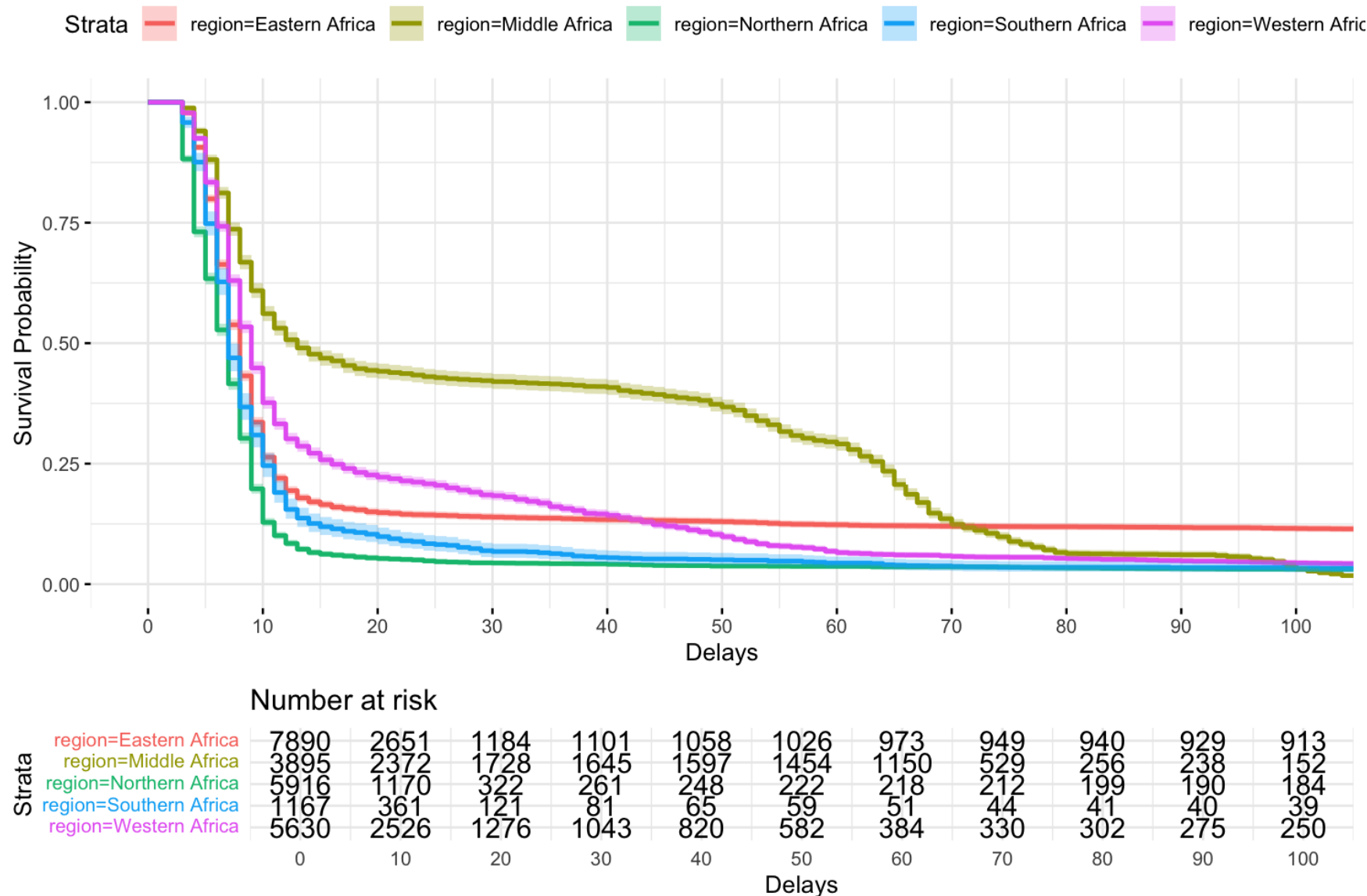This should have been smaller if not for Somalia (East Africa)

# KM (Region)



- East Africa goes down very fast between 0 and 14, and then a steady long tail.

- How fast do these curves descend?

- **The effect of the REGION on the Hazard function is less for higher delay?**
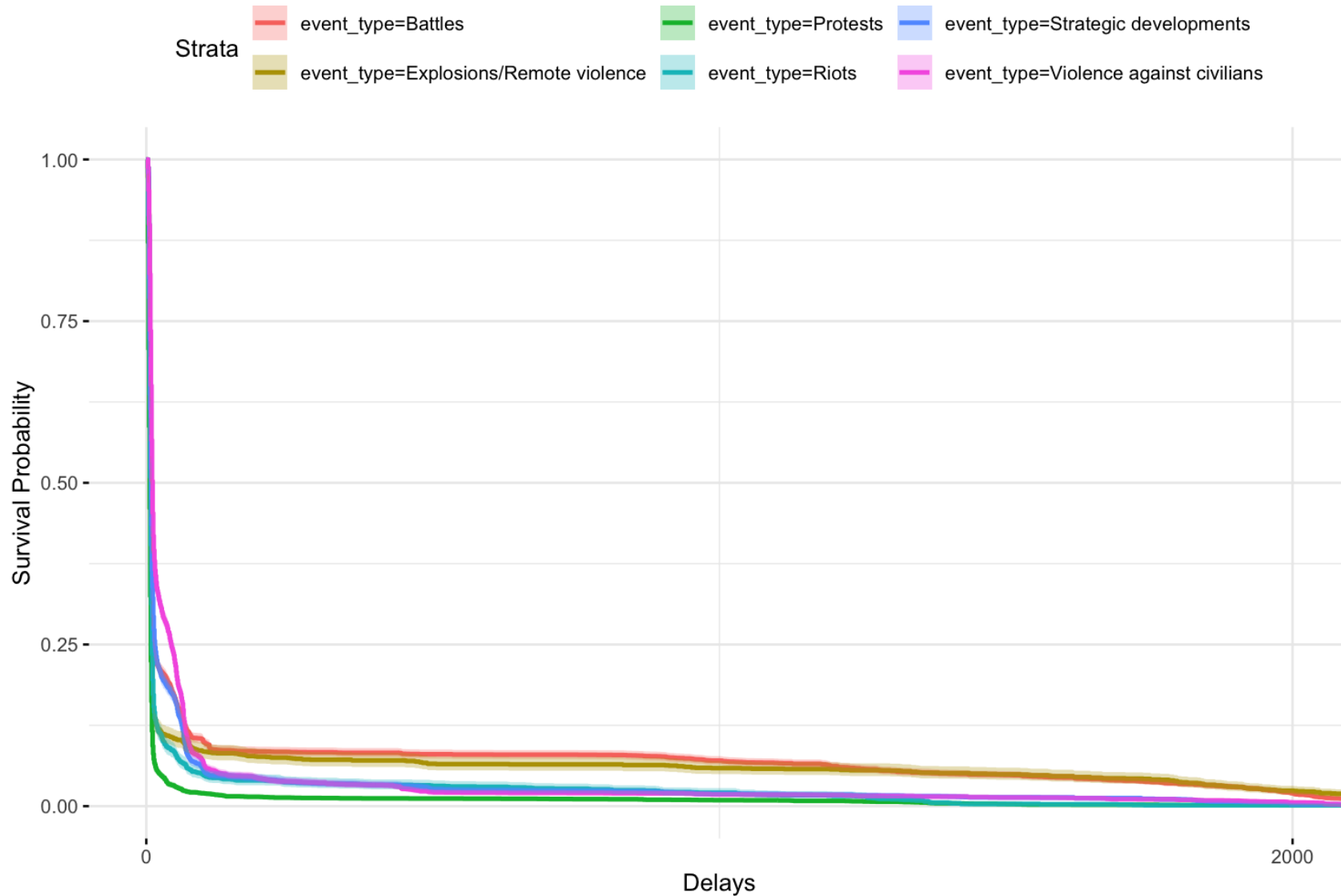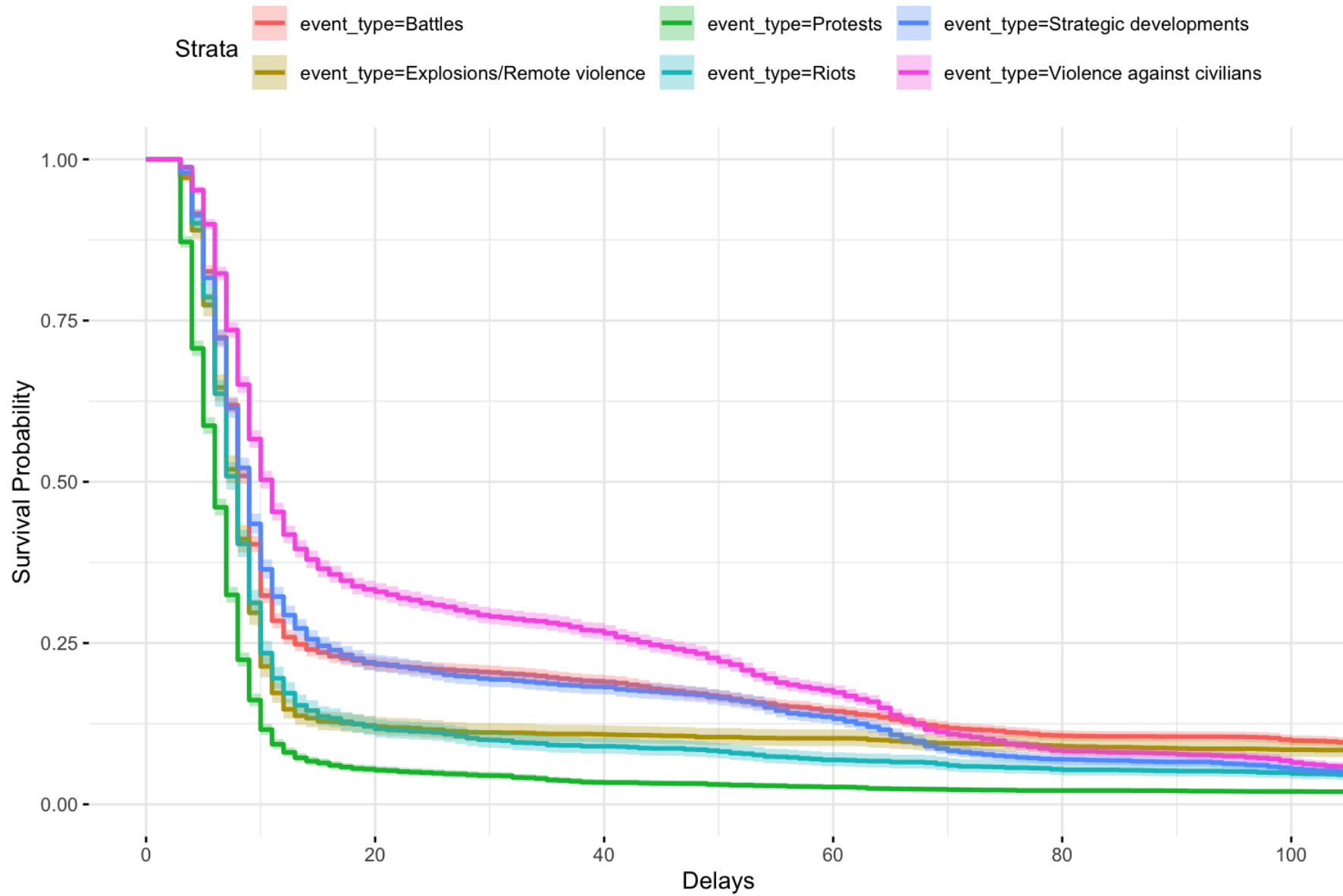
# KM (Region)



- Northern Africa has high number at risk but still decreases very fast

- Middle Africa has the slowest decrease while having the second least number at risk

- East Africa has the highest number at risk and decreases fast before staying steady (because of Somalia)

- **Is there an exponential decay of the number at risk? (Decrease that is function of the amount)**

# KM (Events)



- We observe significant difference between the group (Battles, Explosion) and the rest. We'll also see later that there is also a higher error of event classification in the first group because they are very related and hence more misclassified [battles that are misclassified as explosion and almost all of them in Somalia]

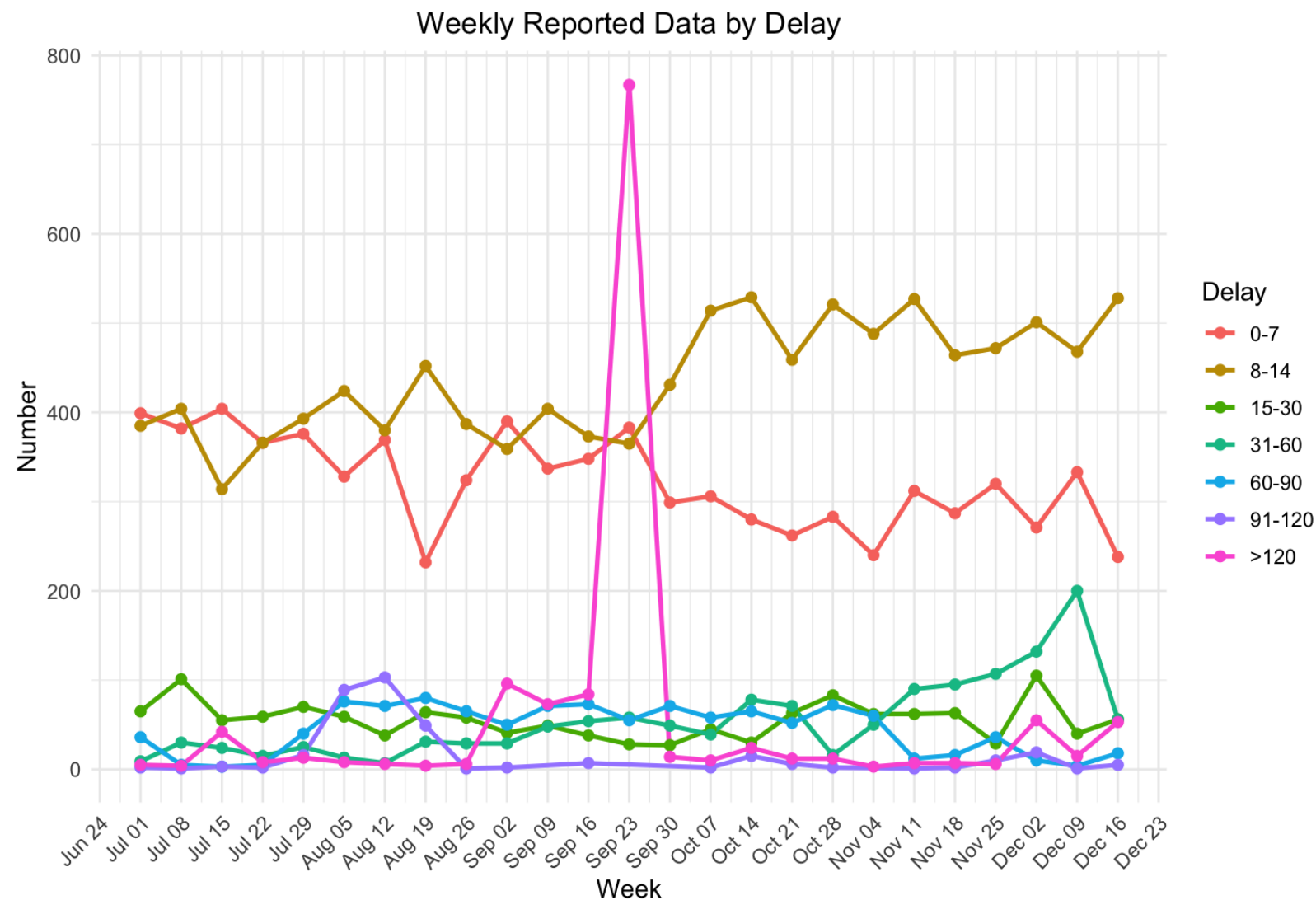# KM (Events)

# What to discuss?

- **Is early truncation then a very smart move here?**

- Unlike the focus on short-term analysis seen in pandemic research, conflict studies require equal attention to both short-term and long-term analyses. This dual focus might be imp for capturing the complex dynamics of conflict over varying timeframes (CYCLE OF CONFLICT vs CYCLE OF PANDEMIC)

- All of the previous results will analytically make more sense when we model the Hazard Function

# Measurement Error effect

- **The measurement error might introduce bias when we estimate the effect of the covariates on delay**

- **How does it change the survival function?**

- **How will we also formalize it into our model? It changes the time at which the true event is reported (possibility: add more events in the class of events with longer delay and removes events from the class of events with shorter delay)**

# Error effect : how do we define it? (TBR)

NO ERROR ASSUMP



Weekly Reported Data by Delay

# Error effect on delays

ANY CHANGE IN ONE OF THE COLUMNS EXCEPT SOURCE

# Error effect on delays

## ANY CHANGE IN ONE OF THE COLUMNS



Weekly Reported Data by Delay

**The events with short delay did not change**

# Error effect on delays

THE LATEST TIMESTAMP



Weekly Reported Data by Delay
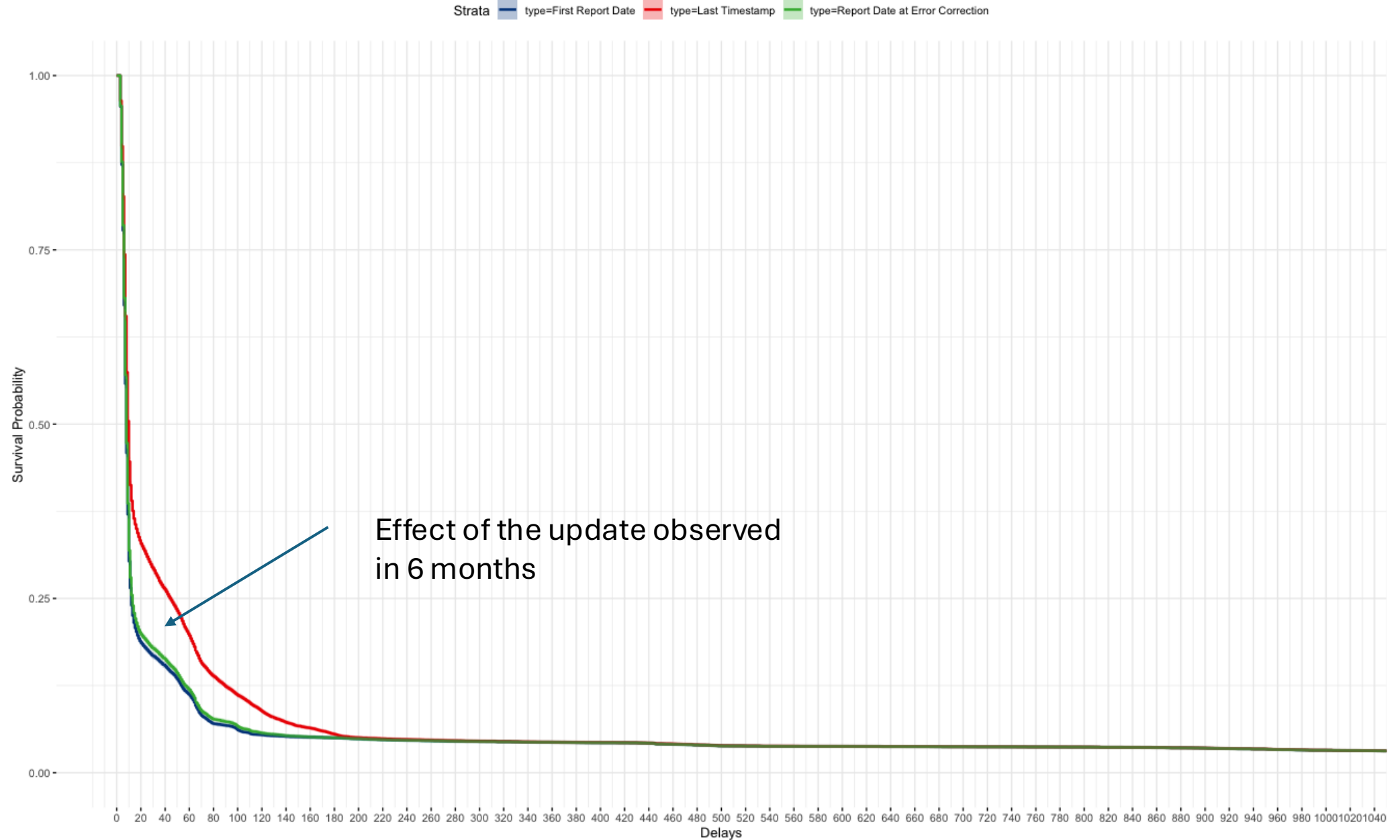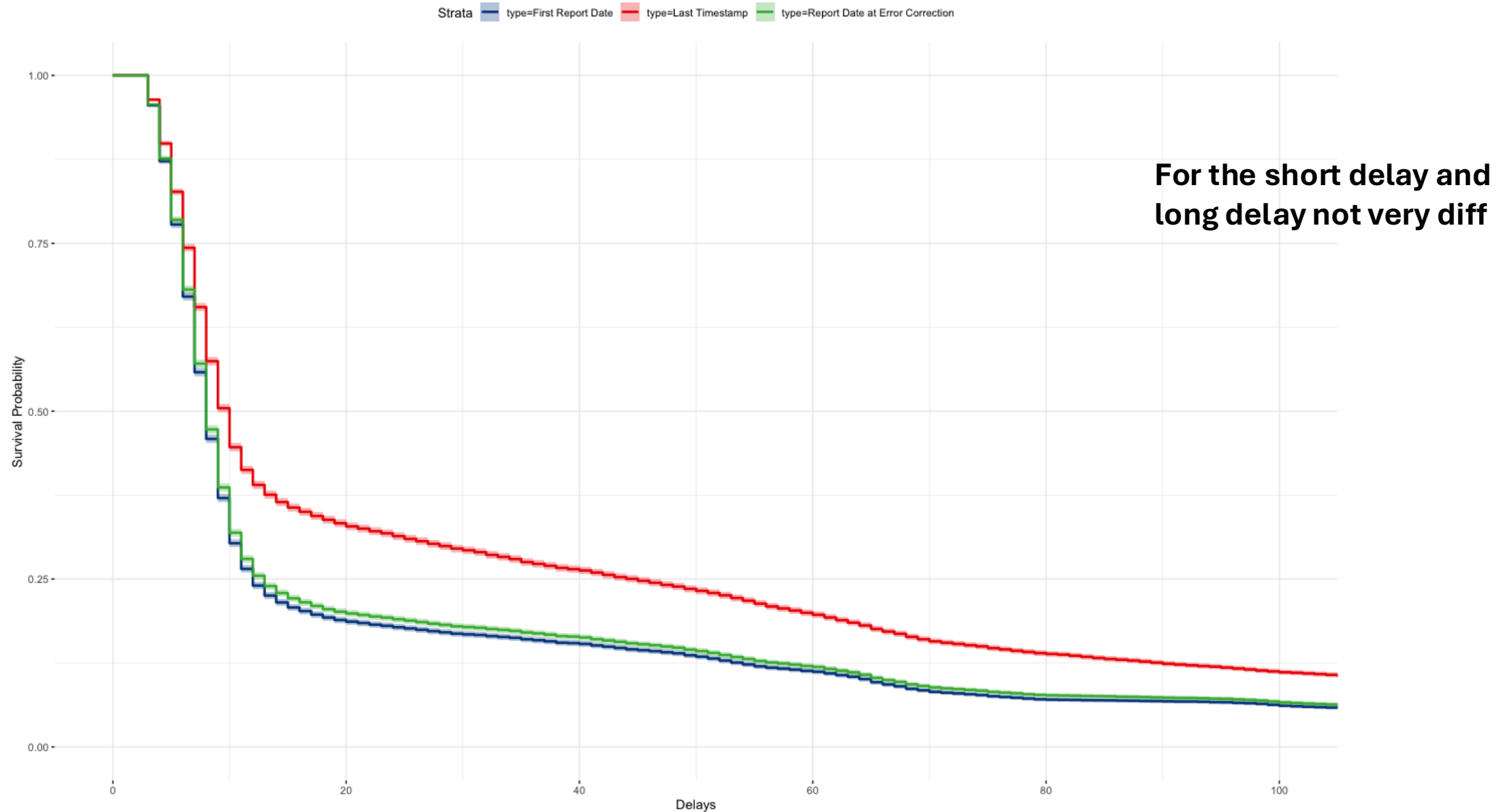
# Error effect on delays

- How our KM will change depend on how we define the error (what kind of update is considered as error correction : change of event for example, change of location that is more than 60km,...??)

- The columns not figured in our data are not relevant?

# Error effect on delays

# Error effect on delays



**Strata** — type=First Report Date   type=Last Timestamp   type=Report Date at Error Correction

**For the short delay and long delay not very diff**

# A bit of discussion

- Are these differences significant? Is it worth modelling? (Error observed on data collected in 6 months, is there a way to foresee how error grows with higher period of observation)

- Will we tackle this hierarchically? (nowcast model then error correction based on it, or error correction then nowcast model based on the corrected data)

- We observe very few error corrections among newly reported events because we observed the change for only 6 months. One extreme case that we would fail to notice with such data is a negative relationship between delay of report and delay of correction from report (**event with long delay gets corrected rapidly and event with short delay gets corrected after a long time** )

- Can we gain knowledge about the error correction from the updates? (are they proportional)

- How long does correction take place?

# Observational Error: Descri

- 164 changes in **date event** observed between first and following downloaded datasets (up to 6 months)
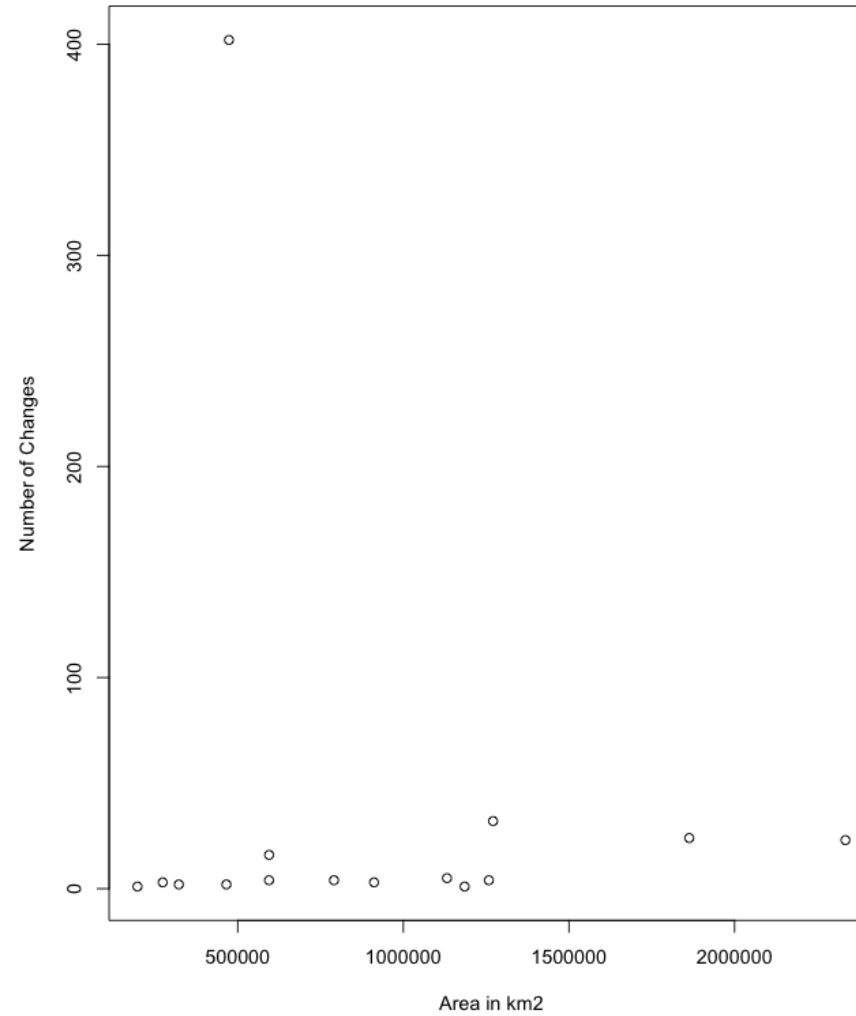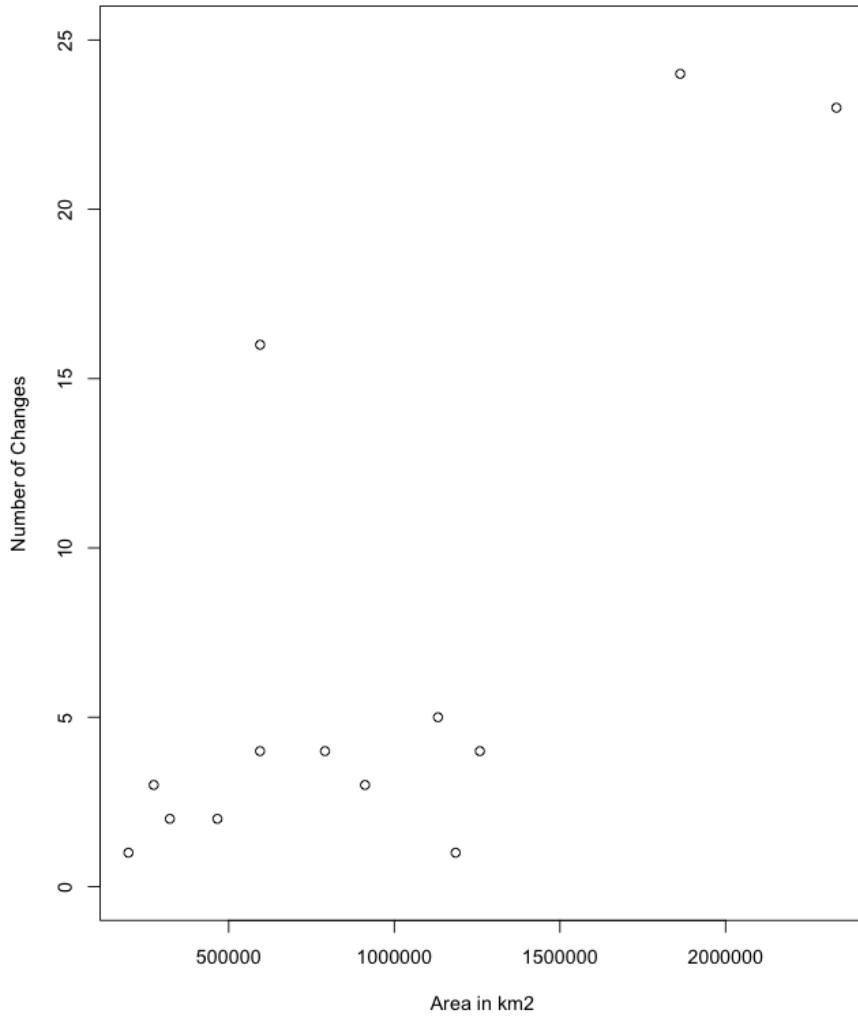
| Diff | 1 | 2 | 3 | 4 | 5 | 7 | 9 | 11 | 15 | 21 | 25 | 31 | 40 | 61 | 274 |
|------|-----|----|---|---|---|---|---|----|----|----|----|----|----|----|-----|
| N | 125 | 13 | 7 | 4 | 1 | 2 | 1 | 1 | 1 | 2 | 1 | 3 | 1 | 1 | 1 |

- 105 changes in **event type** observed in 6 months

| Burkina Faso | | Cameroon | Central African Republic | | Comoros | Democratic Republic of Congo | | Egypt |
|---|---|---|---|---|---|---|---|---|
| 5 | | 3 | 1 | | 1 | 3 | | 1 |
| Ethiopia | | Ghana | Guinea | | Kenya | Madagascar | | Mali |
| 3 | | 5 | 3 | | 10 | 3 | | 3 |
| Morocco | | Mozambique | Niger | | Nigeria | Somalia | | South Africa |
| 1 | | 6 | 1 | | 5 | 30 | | 2 |
| Sudan | | Tanzania | Togo | | Uganda | | | |
| 15 | | 2 | 1 | | 1 | | | |

Interesting fact : The level of corruption in countries with relatively high event changes is also relatively higher

# Observational Error: Descri



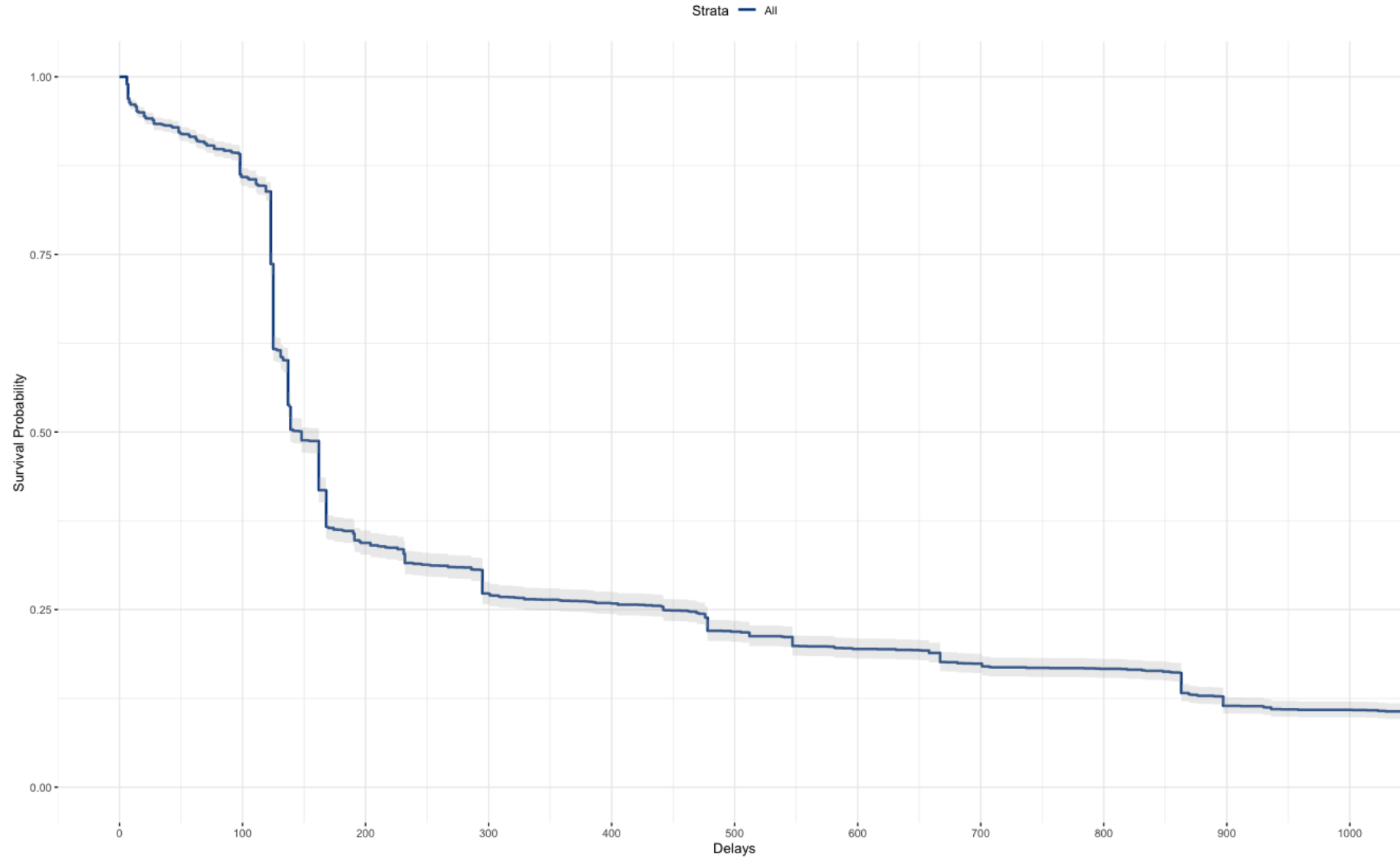Change of location (more
than a radius of 60km)
Somalia

# Observational Error: Descri

| event_type | next_event Battles | Explosions/Remote violence | Protests | Riots | Strategic developments | Violence against civilians |
|---|---|---|---|---|---|---|
| Battles | 0 | 0 | 0 | 0 | 5 | 2 |
| Explosions/Remote violence | 32 | 0 | 0 | 0 | 1 | 1 |
| Protests | 0 | 0 | 0 | 18 | 2 | 0 |
| Riots | 1 | 0 | 6 | 0 | 0 | 2 |
| Strategic developments | 1 | 0 | 0 | 0 | 0 | 7 |
| Violence against civilians | 17 | 3 | 1 | 4 | 2 | 0 |

28 out of the 32 in Somalia

These corrections make sense. Events with higher correction happen to be more correlated

# How long do error get corrected

Is there a pattern we can pick in the error correction (something our model can statistically learn)
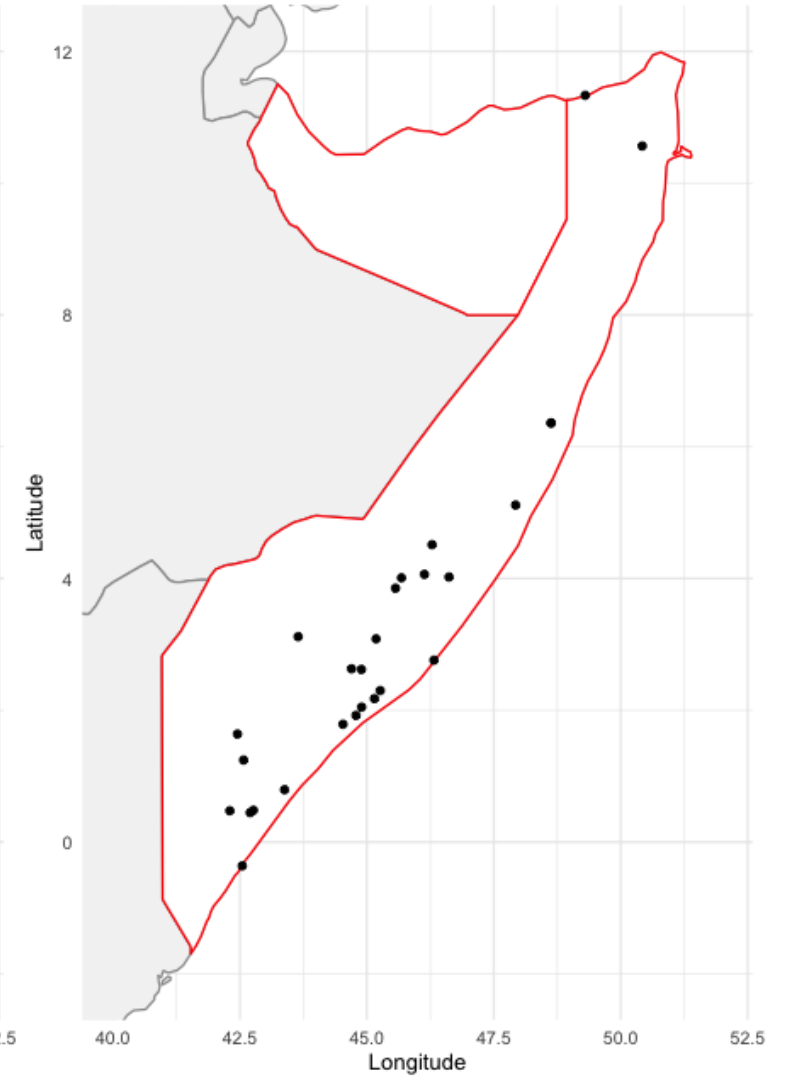
# IN GENERAL

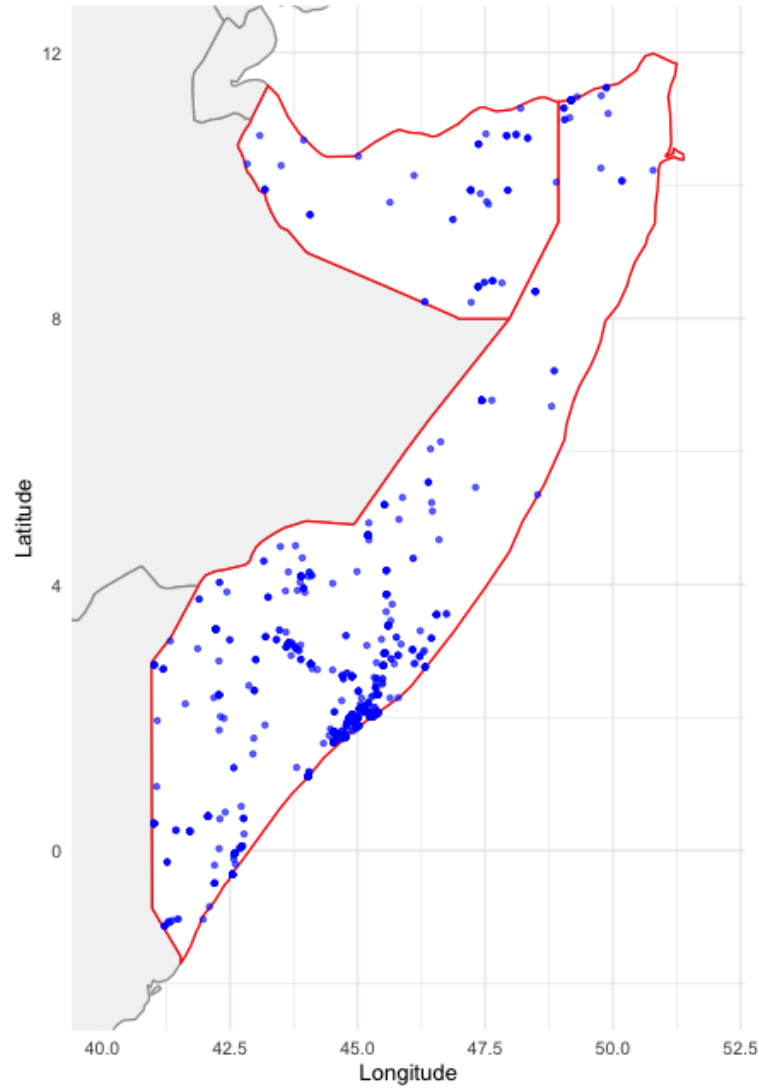

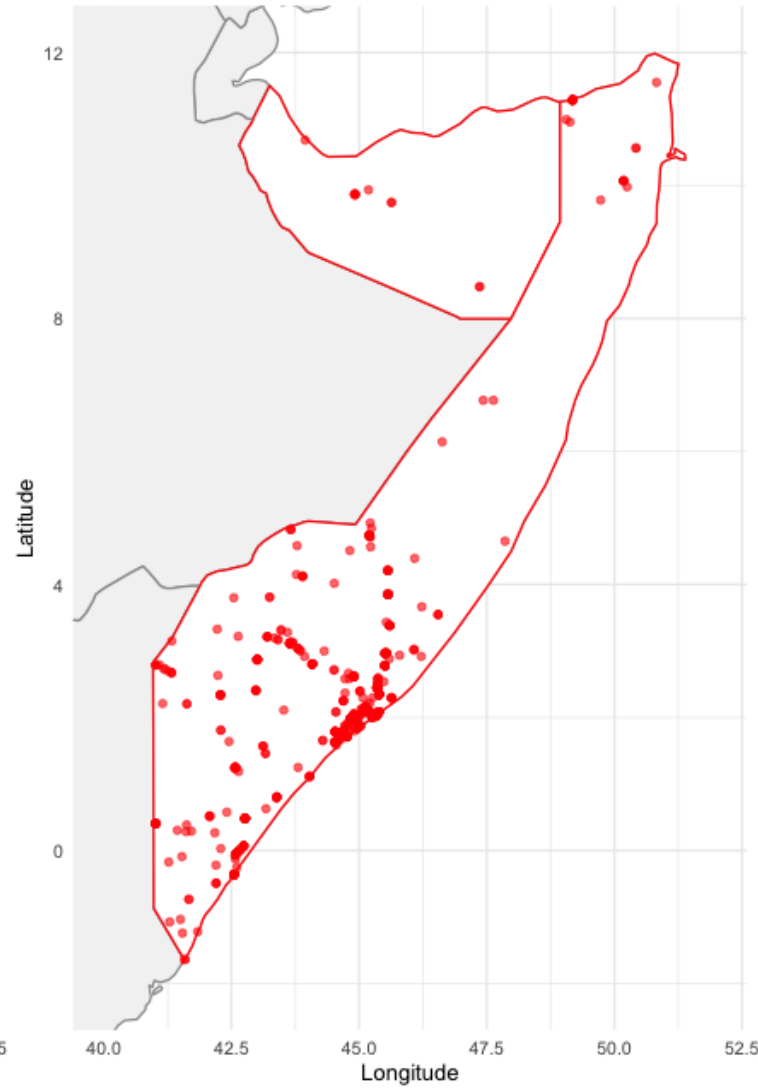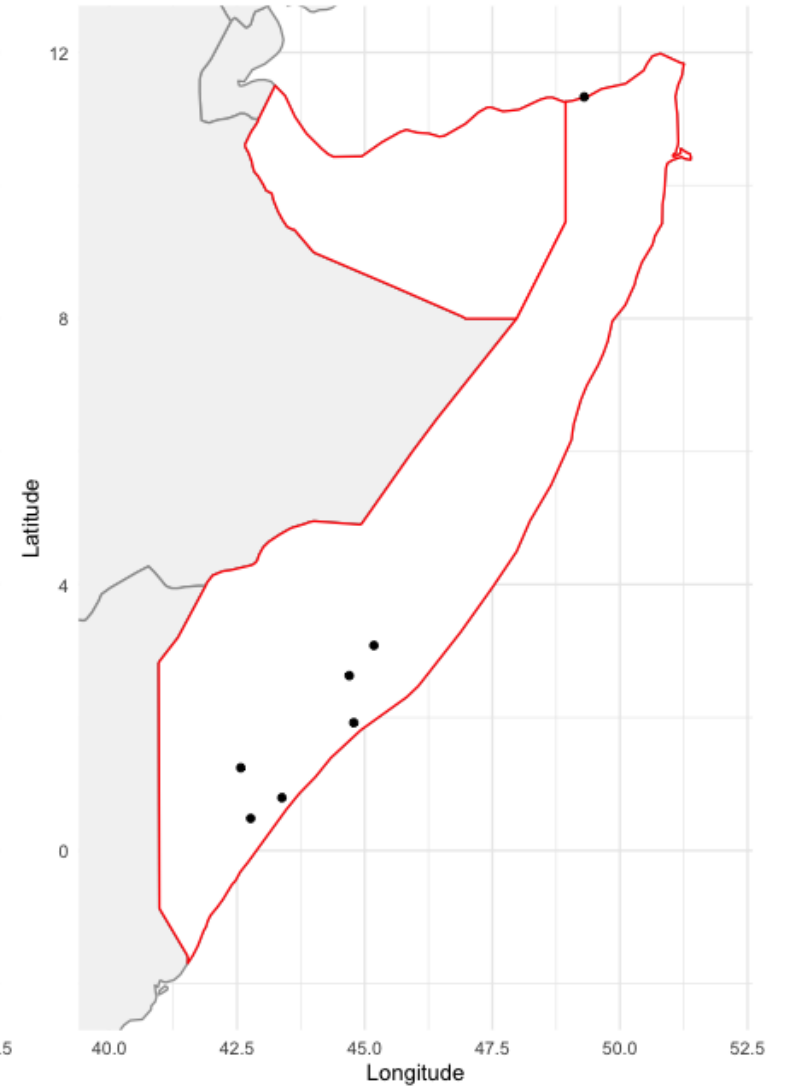Battles in Somalia — Explosion in Somalia — Explo to Batt in Somalia

# 2019

# 2020

# 2021



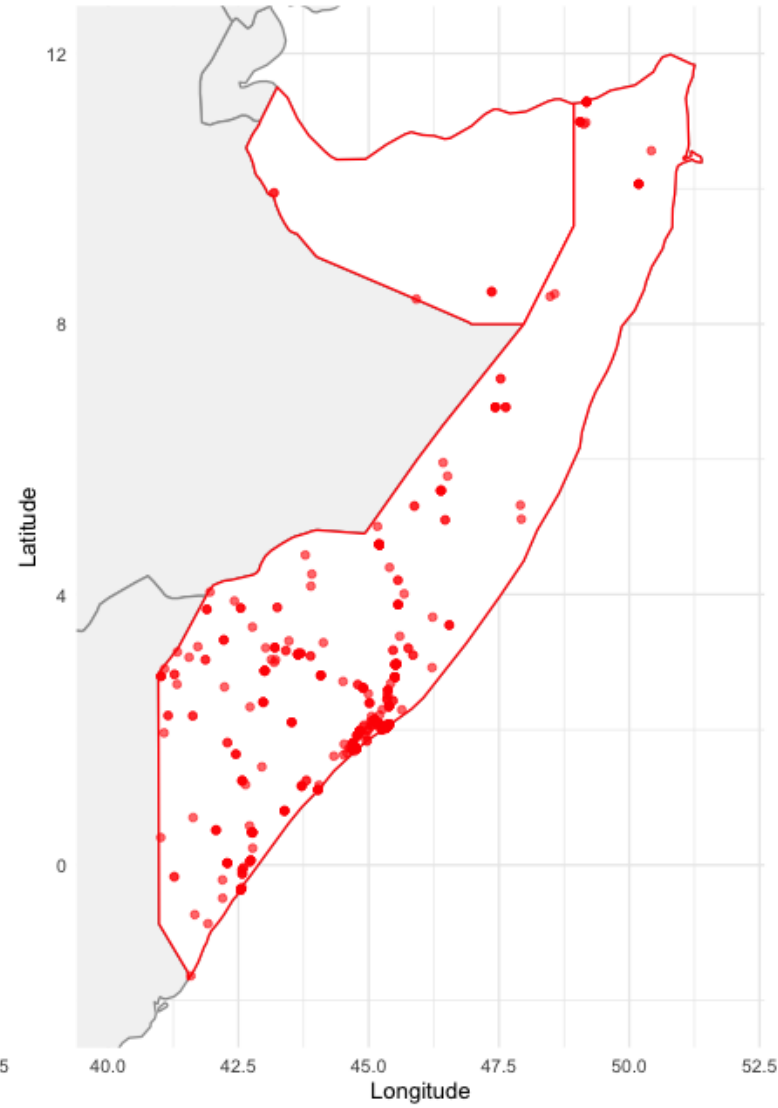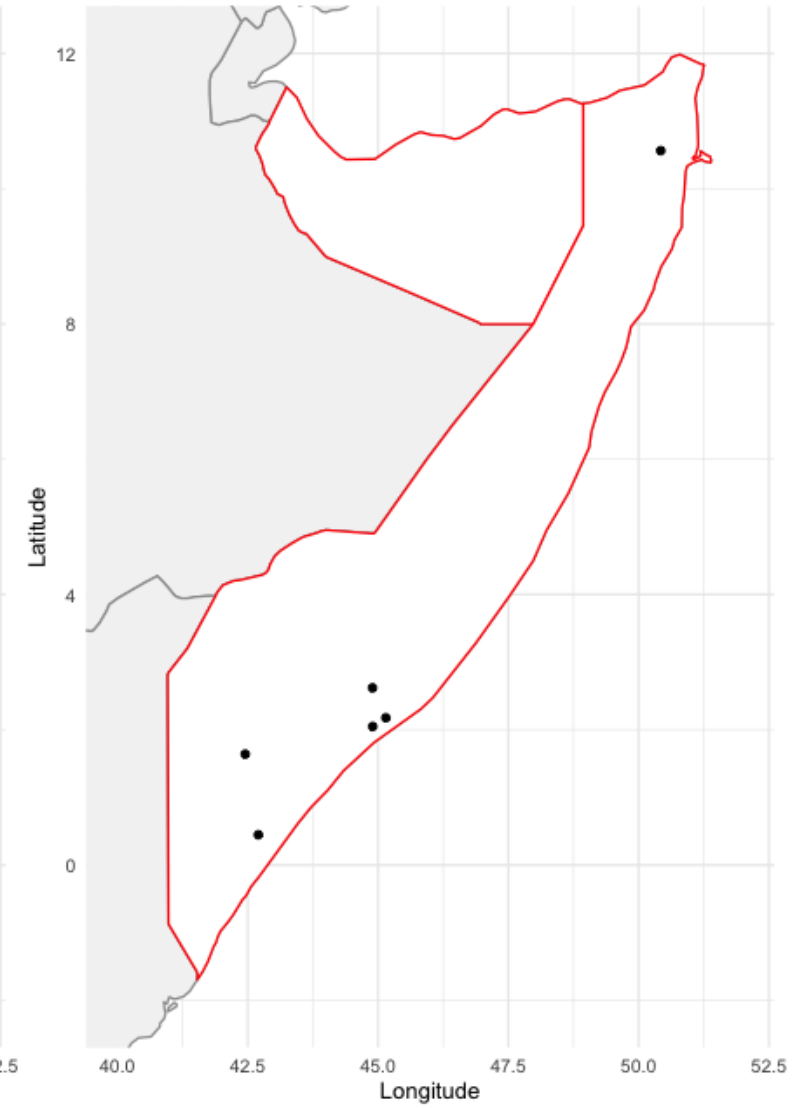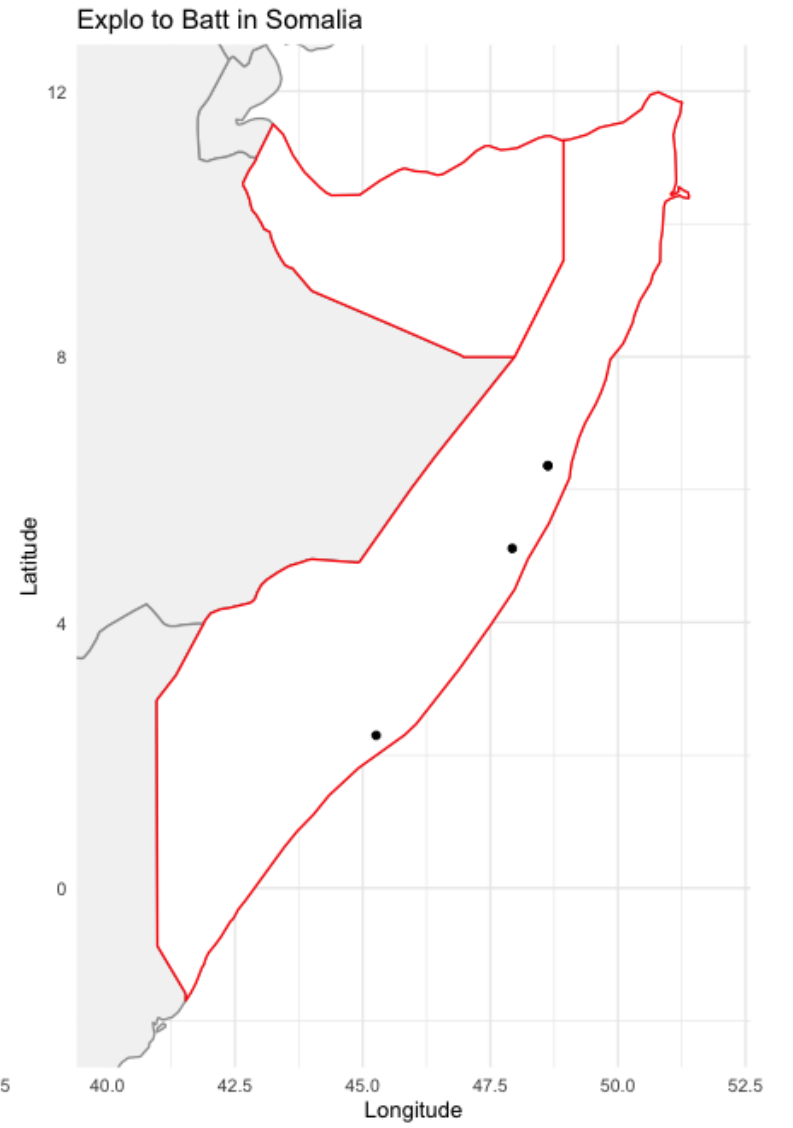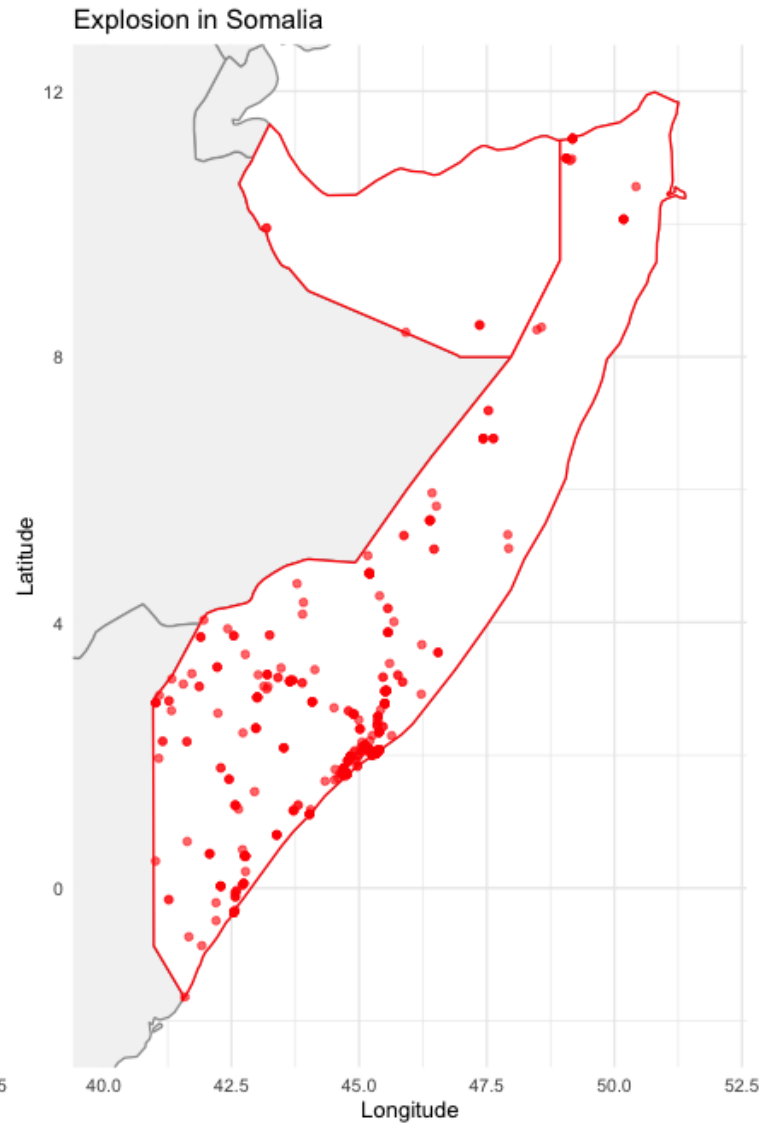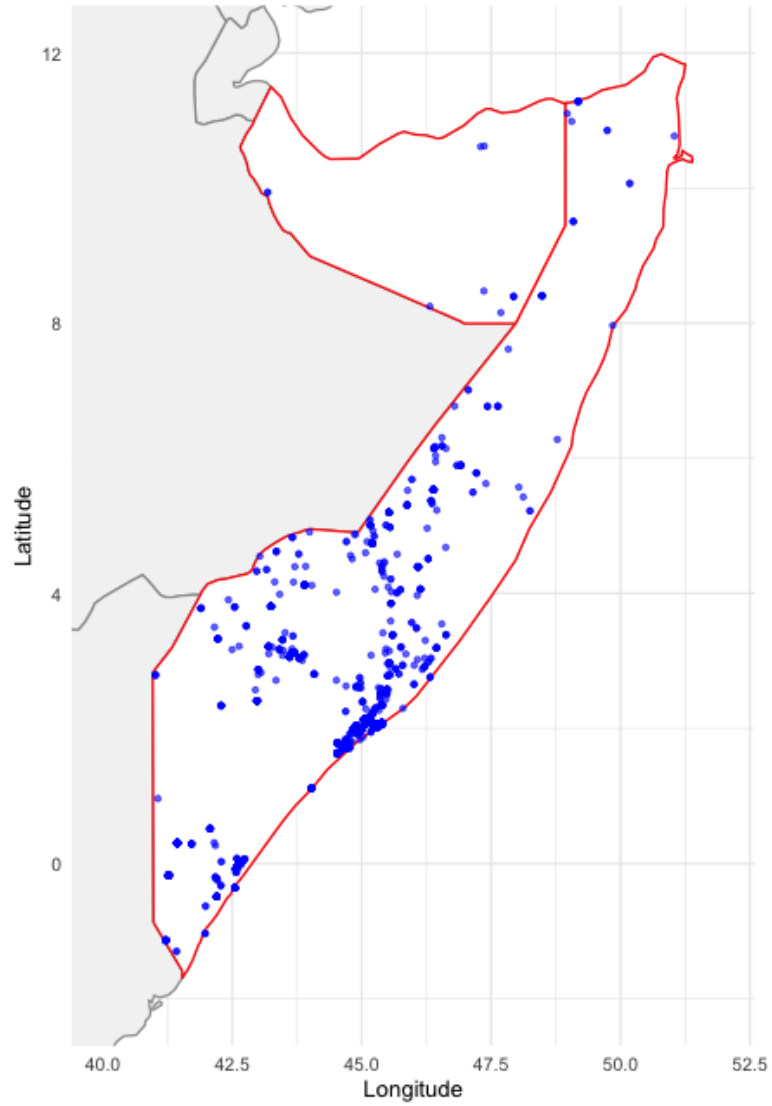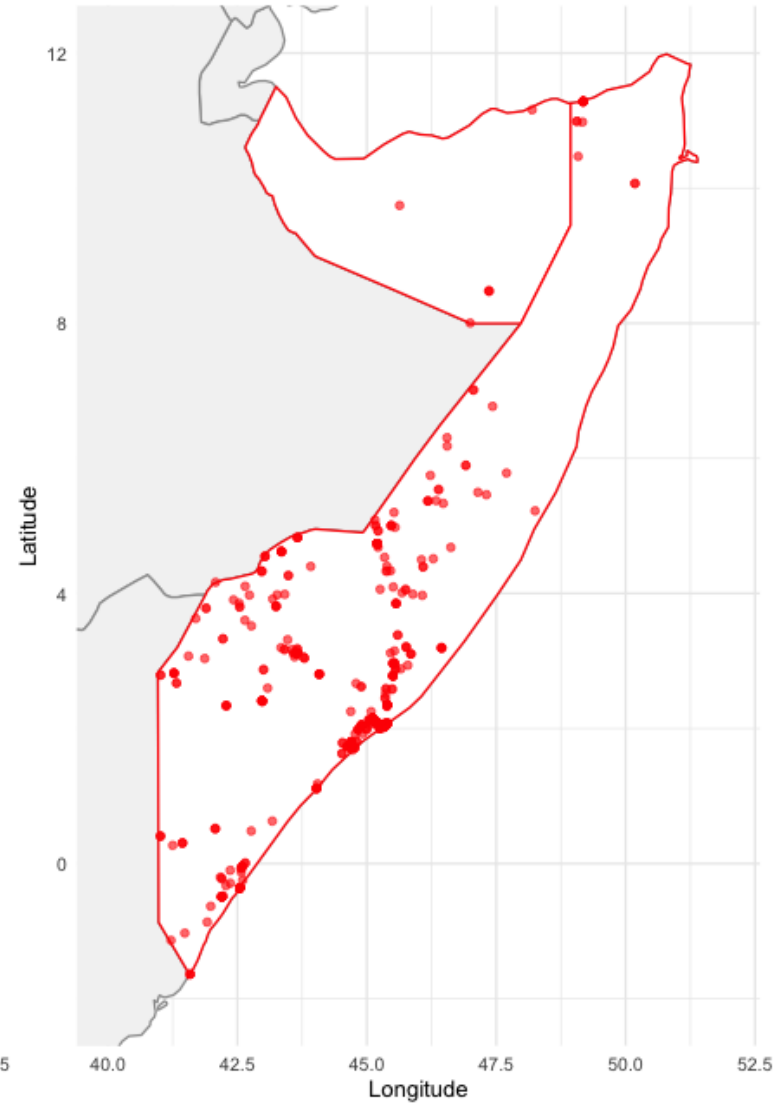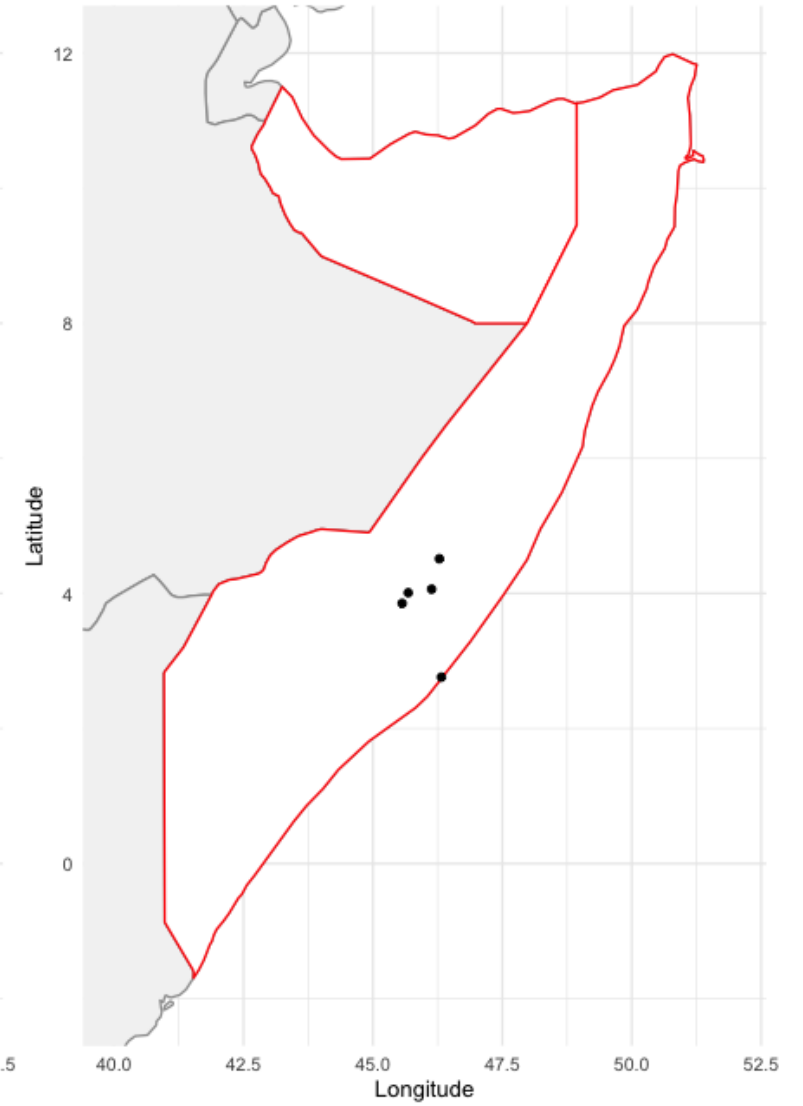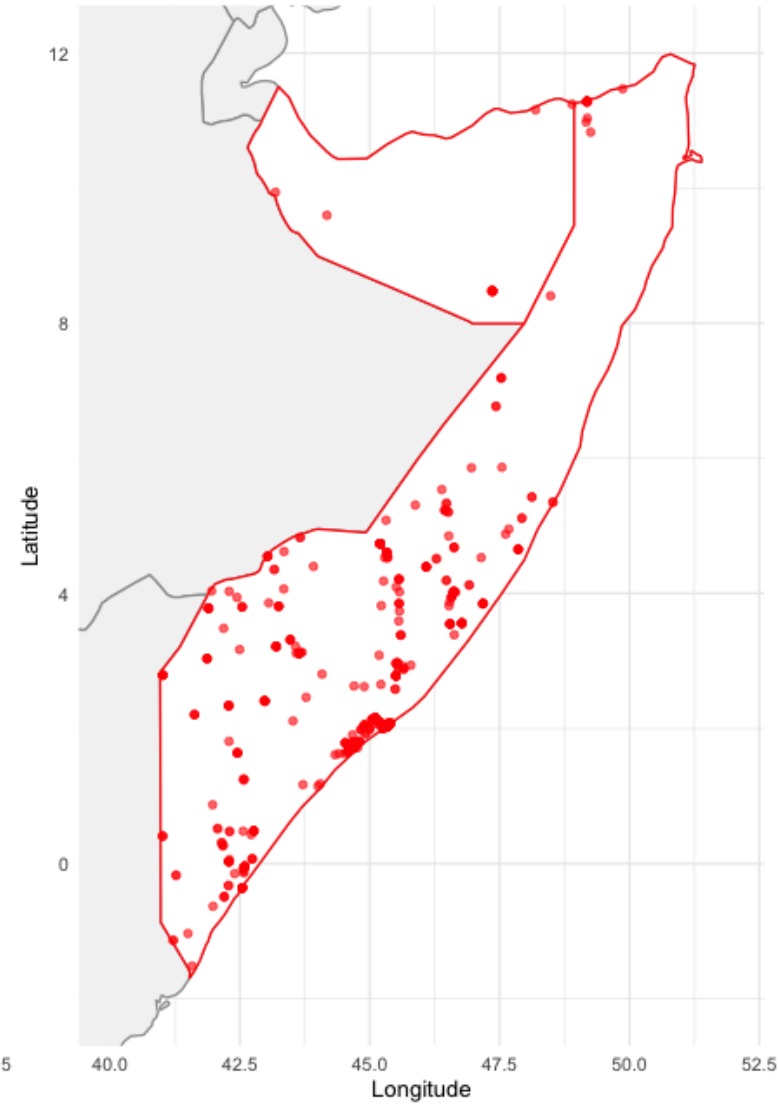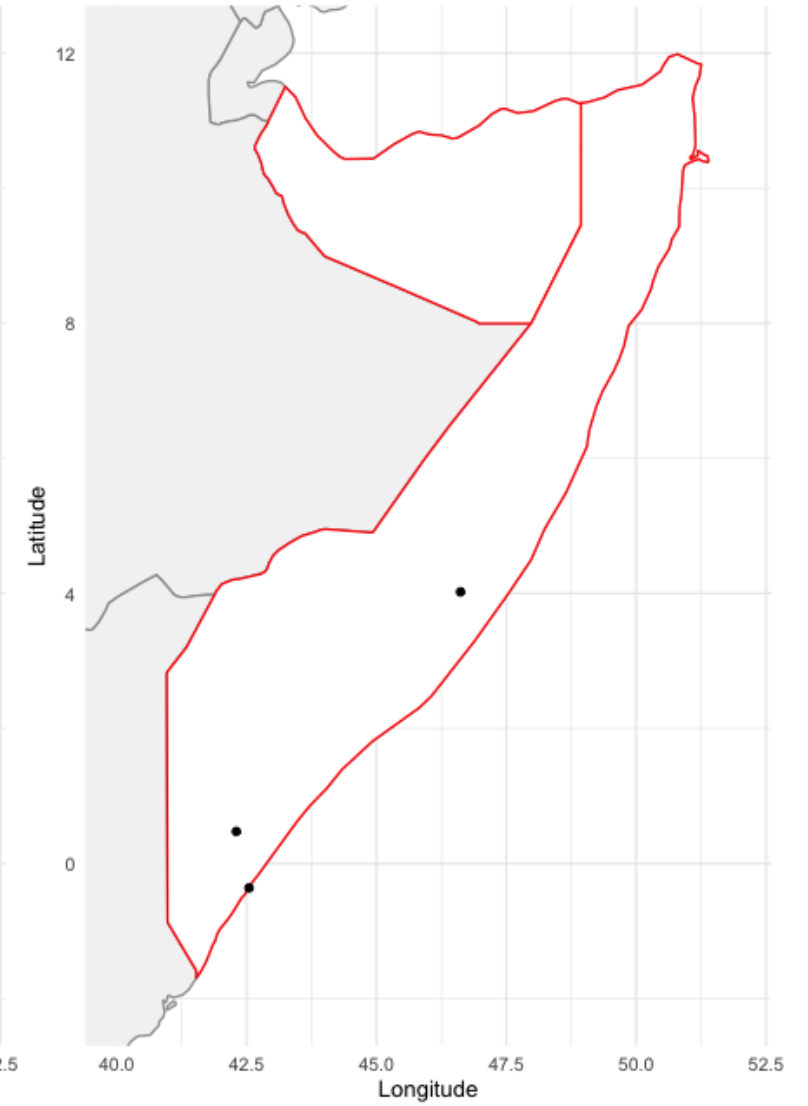Battles in Somalia · Explosion in Somalia · Explo to Batt in Somalia

# 2022

# 2023

# 2024

- What we can possibly assume from this :

    - There is an interaction between event types (intra), and between events and location. There are some preserve patterns over the years
    - How the "reported" battles and explosion change over the years are proportional. This is seen from how their KM overlaps (even the long tail)
    - Correction does not occur in the short term (we will analyse the delay of correction more thoroughly later). Older events (wrt time of event, not report time) get more chance to be corrected. Maybe the 2024 reported events that we know of now will be corrected in 4 years (something we need to foresee)
    - The older events had more chance to get corrected throughout the years, yet we still observe some correction in 5 months. With more months of observation, we could observe more error correction (the measurement error is then underestimated if we use only the correction during the 5 months)

- **What we do not know.**
  - o The time of report of every event. Which leaves us in doubt because it is also possible that delay of correction is negatively correlated to the delay of report. Which is why we do not have any correction for reported events that happened in 2024 (because they have less than a year of delay report: that at least we know)
  - o The error. We only observed the correction in 5 months. But is it enough to say anything about the measurement error given our data exhibit a long tail survival function for delay?

# Updates Behaviour

- Can we deduce any useful information from the timestamp of events reported on or before the first download AND did not get updated anymore?

- Problem : The timestamp is ambiguous. It can signify
  o Time of first report
  o Time of last update without error correction
  o Time of last update with error correction

$$T_s = T_e + D_r + \sum_{i=1}^{p} U_i$$

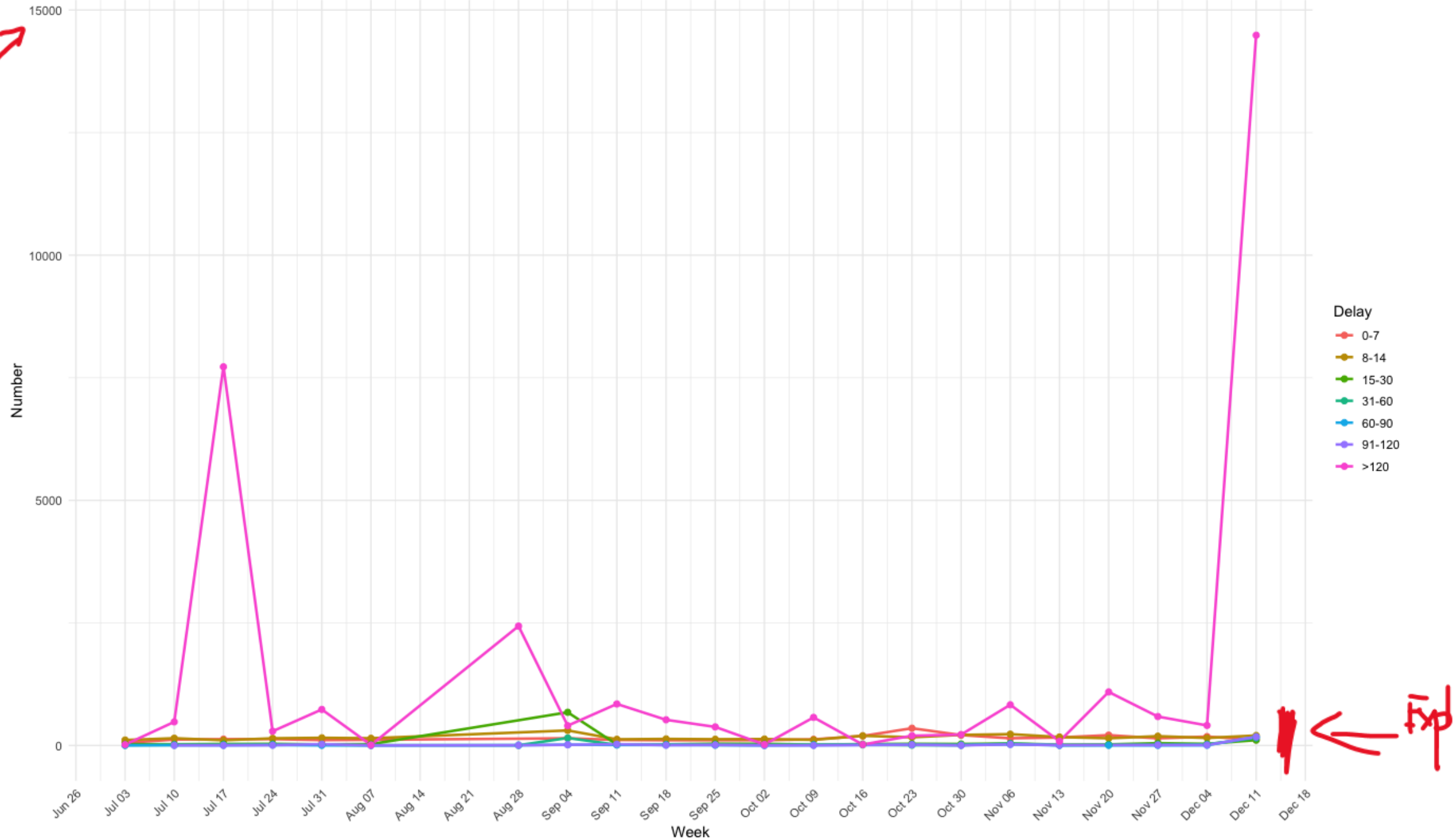Ts : Last timestamp
Te : Time of event
Dr : Delay of first report
Ui : Delay of i-th update after the (i-1)th
where $U_0$ = Dr
p : number of updates after first report

- **Can we then assume that the distribution of Dr and Dr+Sum(U) is the same? (negligeable error correction)** <span style="color:red">NO</span>

- Assuming delay distribution is **stationary,** the KM delays of the events newly reported between our first dowload and the followings AND the KM delays of the events reported during the same period the year before by using the difference between time of event and timestamp **should not have significant difference** (just with errors due to randomness) <span style="color:red">WRONG</span>

Weekly Reported Data by Delay

# Weekly Reported Data by Delay



**Delay**
- 0-7
- 8-14
- 15-30
- 31-60
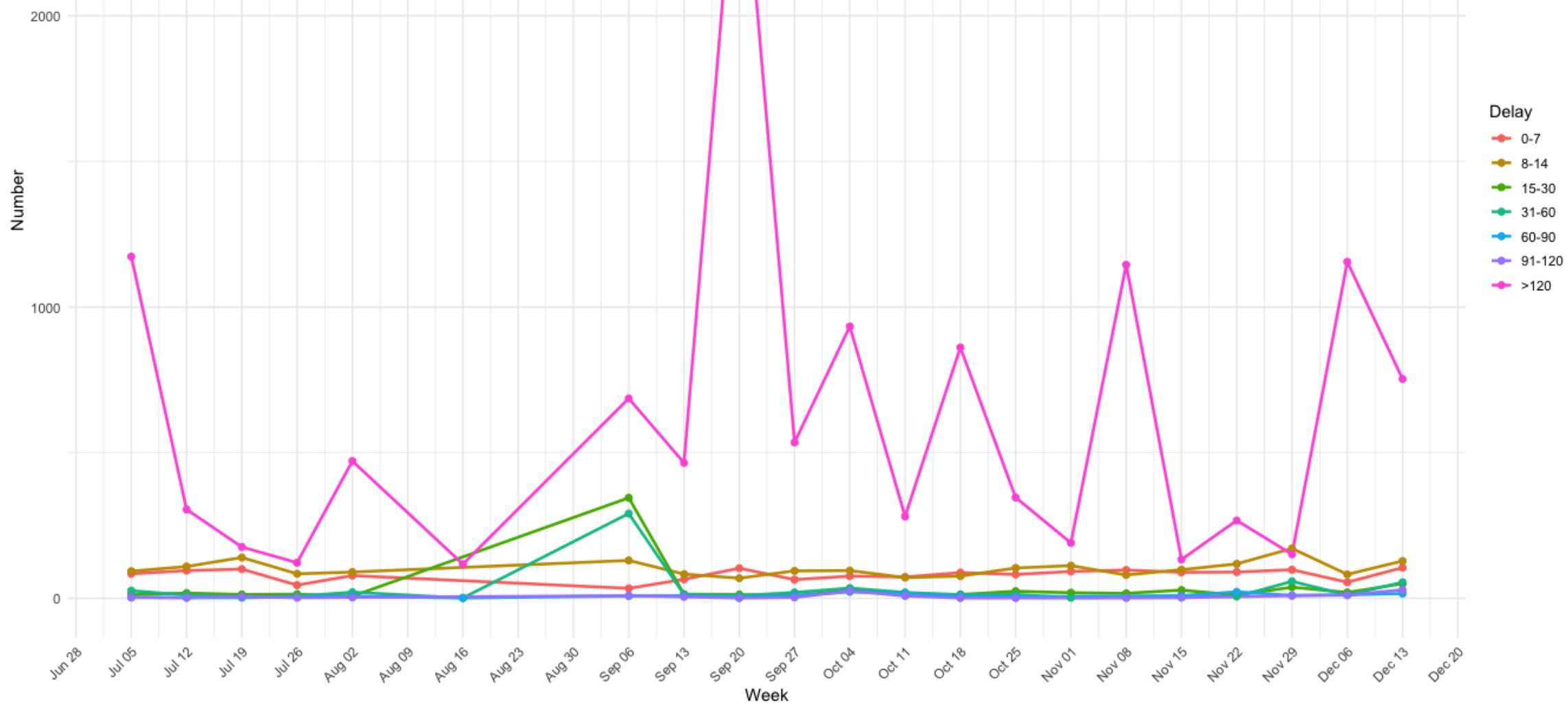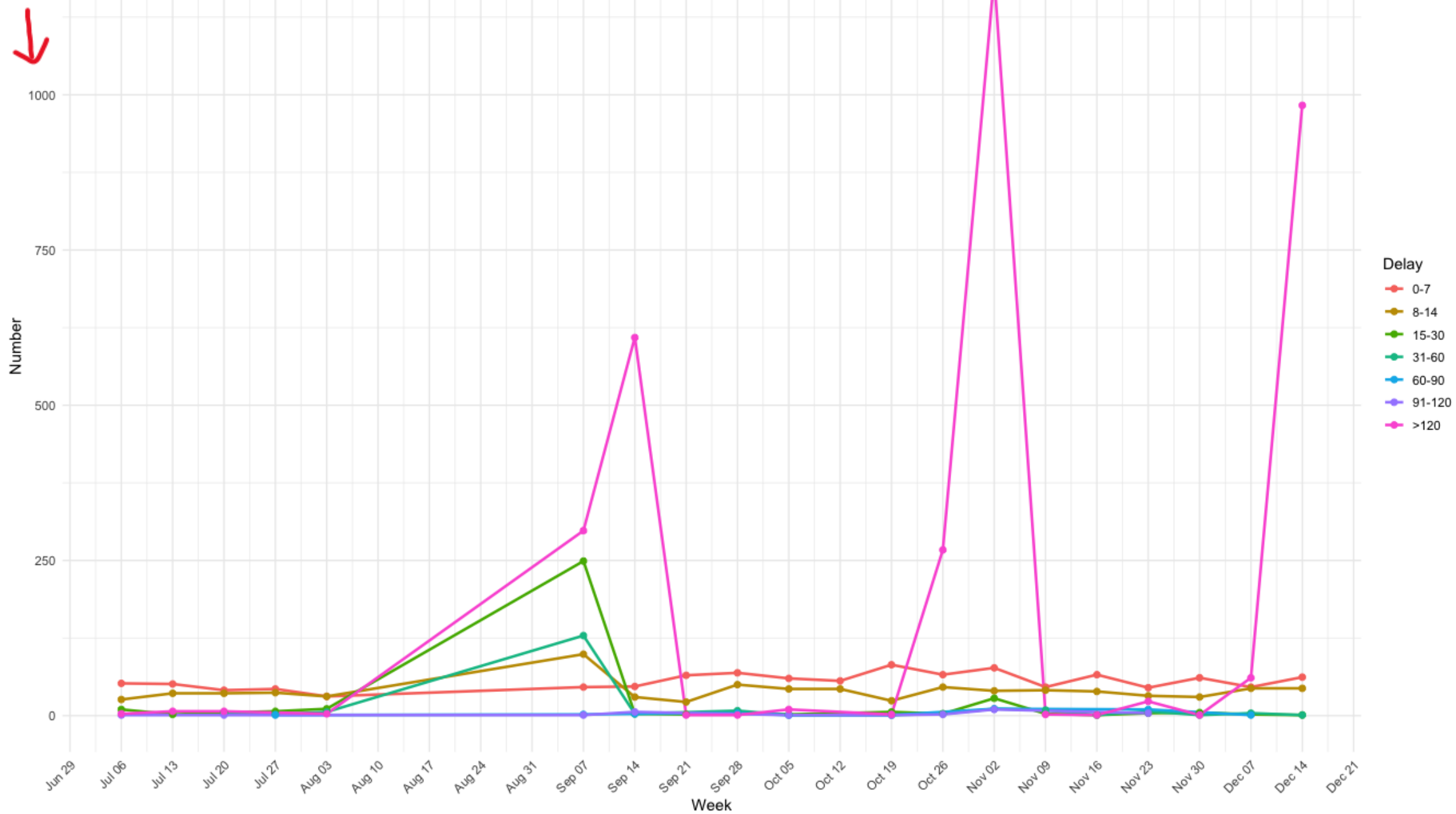- 60-90
- 91-120
- >120

Number

Week

Weekly Reported Data by Delay
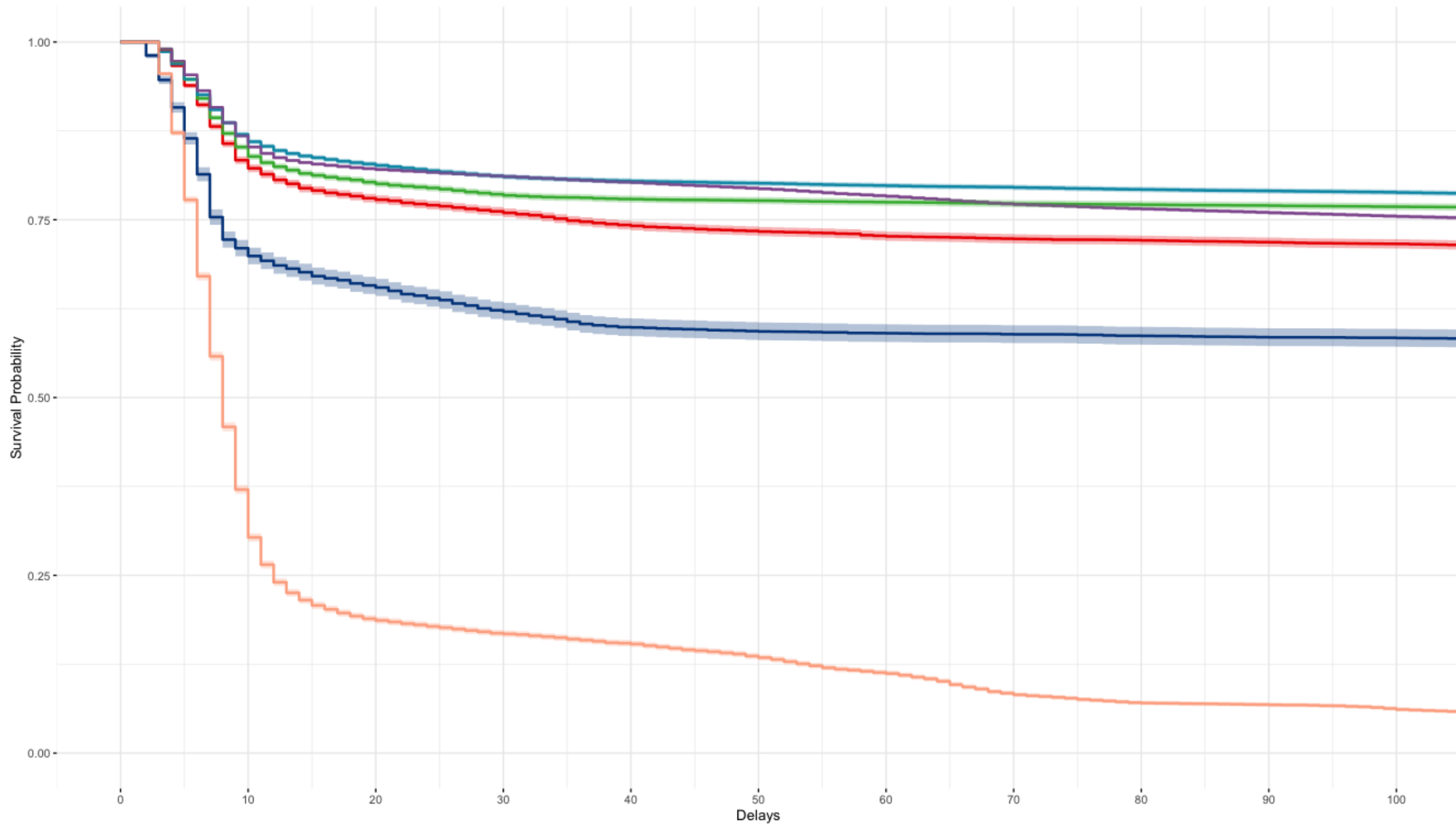
Weekly Reported Data by Delay

# Discussion

- It's either we have consistently relatively high number of late reports or late updates or both (compared to the event time)

- The number of events with less delay of report and cumulated delay of update is very small. These indicate the events that give us more evidence that they didn't get updated (because Dr+Sum(U) is small). The inverse is not true. The events with high Dr+Sum(U) could be events with high Dr or high Sum(U) or both

- The number of events reported/updated decreases every year. A rational explanation is that there is a migration from preceding to the next years. Hence we have many updates. And the older the event is the higher is the probability of it being updated (or corrected)

$$S_{2023}(t) \geq S_{2022}(t) \geq S_{2021}(t) \geq S_{2020}(t)$$

In the short delay

We see big difference between KP 2024 (where we took only new events and no update) and the others where we took the last date of update.

**Hence the update is then very important because many of the reported events will be updated.**

# Another discussion

- We talked about update but not every update is important enough to be considered "error". However we could prove with common sense that update is underestimated using only updates observed during the 5 months. The latter increases the evidence that our error is also underestimated (with the assumption that error is proportional to the updated).

- How can we then estimate the "relevant" part of the update. Our best shot is to analyse these changes that occurred in the 5 months of observation and find a pattern.

- If we have date of event, covariates, last update, can we tell with some level of confidence if it was an event that got corrected? And got corrected from what (ex : Explo to Battle?) and when did it get corrected?

A few late reports in general. Except the massive report concentrated on some date points

Exploratory

High difference between event and timestamp implies very high odds for updated event (except the concentrated ones, because this patter happened even with the 2024

Can we fit a model that learn this

Which updates are more likely to be an error correction

Error correction with respect to event or location and how long did it take to be corrected.

We know error behavior