

# ZAi Notes

Lukas Bayer

November 17, 2023

# 1 Fake-Factors

## 1.1 Nomenclature

We use the following samples to estimate the amount of fake electrons in our data:

`data` is a sample of measured data from the ATLAS experiment. It includes real as well as fake electrons.

`MJ` is the subset of the data sample that consists only of fake electrons from multi-jet events. A priori, its size and distribution are unknown.

`prompt` is a sample of Monte-Carlo generated events that consists only of simulated real electrons. It is normalized to the luminosity, such that, for a perfect simulation, it would include the same number of real electrons as the data sample. In reality, they differ by a factor of  $\hat{\mu}$ , which is to be determined.

$$N_{\text{MJ}} = N_{\text{data}} - \hat{\mu} \cdot N_{\text{prompt}} \quad (1)$$

Additionally we define the identification regions:

`ID` are events passing tight identification criteria

`nL` are events failing loose identification criteria<sup>1</sup>

and invariant mass regions:

`on` are events with an invariant mass close to the Z pole  $m_{\ell\ell} \approx m_Z$

`off` are events with an invariant mass far away from the Z pole  $m_{\ell\ell} \gg m_Z$

thus dividing each sample into a total of four regions. `SR = IDon` is the signal region of interest, while the remaining three regions `nLon`, `nLoff` and `IDoff` are enriched in multi-jet experiments and used as control regions `CR`.

## 1.2 Calculation of $\hat{\mu}$

Assume that the ratio of fake electrons passing to failing identification is the same on and off the Z pole. Then the amount of fake electrons in the signal region is given by

$$N_{\text{MJ}}^{\text{IDon}} = N_{\text{MJ}}^{\text{nLon}} \cdot \frac{N_{\text{MJ}}^{\text{IDoff}}}{N_{\text{MJ}}^{\text{nLoff}}} \quad (2)$$

Rearranging equation 2 and inserting equation 1 leads to:

$$\begin{aligned} 0 &= N_{\text{MJ}}^{\text{nLon}} \cdot N_{\text{MJ}}^{\text{IDoff}} - N_{\text{MJ}}^{\text{IDon}} \cdot N_{\text{MJ}}^{\text{nLoff}} \\ &= \left( N_{\text{data}}^{\text{nLon}} - \hat{\mu} N_{\text{prompt}}^{\text{nLon}} \right) \cdot \left( N_{\text{data}}^{\text{IDoff}} - \hat{\mu} N_{\text{prompt}}^{\text{IDoff}} \right) - \left( N_{\text{data}}^{\text{IDon}} - \hat{\mu} N_{\text{prompt}}^{\text{IDon}} \right) \cdot \left( N_{\text{data}}^{\text{nLoff}} - \hat{\mu} N_{\text{prompt}}^{\text{nLoff}} \right) \\ &= A\hat{\mu}^2 + B\hat{\mu} + C \end{aligned}$$

---

<sup>1</sup>the identification criteria employed here are not necessarily the same as for the scale-factor calculation

with

$$\begin{aligned}
A &= N_{\text{prompt}}^{\text{nLon}} \cdot N_{\text{prompt}}^{\text{IDoff}} - N_{\text{prompt}}^{\text{IDon}} \cdot N_{\text{prompt}}^{\text{nLoff}} \\
B &= N_{\text{data}}^{\text{IDon}} \cdot N_{\text{prompt}}^{\text{nLoff}} - N_{\text{data}}^{\text{nLon}} \cdot N_{\text{prompt}}^{\text{IDoff}} + N_{\text{prompt}}^{\text{IDon}} \cdot N_{\text{data}}^{\text{nLoff}} - N_{\text{prompt}}^{\text{nLon}} \cdot N_{\text{data}}^{\text{IDoff}} \\
C &= N_{\text{data}}^{\text{nLon}} \cdot N_{\text{data}}^{\text{IDoff}} - N_{\text{data}}^{\text{IDon}} \cdot N_{\text{data}}^{\text{nLoff}}
\end{aligned}$$

Solving for (real, positive values of)  $\hat{\mu}$  gives<sup>23</sup>:

$$\hat{\mu} = \frac{-B + \sqrt{B^2 - 4AC}}{2A} \quad (3)$$

### 1.3 Estimation of Multi-Jet Content

After calculating  $\hat{\mu}$  with equation 3, it can be used to determine estimates for the multi-jet contents of the control regions using equation 1. Finally, these are inserted back into equation 2 to estimate the multi-jet content of the signal region.

The main contribution to the statistical uncertainty is expected to come from the data counts (not from  $\hat{\mu}$  or prompt counts). Hence, for the control regions:

$$\Delta N_{\text{MJ}}^{\text{CR}} \approx \sqrt{N_{\text{data}}^{\text{CR}}}$$

and for the signal region:

$$\begin{aligned}
\Delta N_{\text{MJ}}^{\text{IDon}} &= \sqrt{\sum_{\text{CR}} \left( \Delta N_{\text{MJ}}^{\text{CR}} \cdot \frac{\partial N_{\text{MJ}}^{\text{IDon}}}{\partial N_{\text{MJ}}^{\text{CR}}} \right)^2} \\
&\approx \sqrt{N_{\text{data}}^{\text{nLon}} \cdot \left( \frac{N_{\text{MJ}}^{\text{IDoff}}}{N_{\text{MJ}}^{\text{nLoff}}} \right)^2 + N_{\text{data}}^{\text{IDoff}} \cdot \left( \frac{N_{\text{MJ}}^{\text{nLon}}}{N_{\text{MJ}}^{\text{nLoff}}} \right)^2 + N_{\text{data}}^{\text{nLoff}} \cdot \left( \frac{N_{\text{MJ}}^{\text{nLon}} N_{\text{MJ}}^{\text{IDoff}}}{(N_{\text{MJ}}^{\text{nLoff}})^2} \right)^2}
\end{aligned}$$

### 1.4 Calculation of Fake-Factors

Let's define the fake efficiency as the probability of a fake electron passing identification criteria:

$$\epsilon_f := \frac{N_{\text{MJ}}^{\text{ID}}}{N_{\text{MJ}}}$$

According to [1] the fake-factor is defined as:

$$F := \frac{\epsilon_f}{1 - \epsilon_f} = \frac{N_{\text{MJ}}^{\text{ID}}}{N_{\text{MJ}}^{\text{ID}}}$$

But in our implementation it is instead calculated as:

$$F = \frac{N_{\text{MJ}}^{\text{IDoff}}}{N_{\text{MJ}}^{\text{nLoff}}}$$

Why is this different?

---

<sup>2</sup>In the highly improbable case that  $A = 0$  the result becomes  $\hat{\mu} = -C/A$  instead.

<sup>3</sup>Is this sqrt really always positive?

## 1.5 applyFakeFactor

So far I have described how fake-factors are calculated when the makeFakeFactor flag is set. However, I don't understand, what is being done when applyFakeFactor is set. Can anyone explain?

## References

- [1] ATLAS Collaboration. Tools for estimating fake/non-prompt lepton backgrounds with the atlas detector at the lhc, 2022.