# Measurement of the Photon Identification Efficiency with the Matrix Method Using Neural Networks

Can Süslü

Physikalisches Institut,
University of Bonn, Germany

*can.suslu@ug.bilkent.edu.tr*

October 7, 2022

# Outline

UNIVERSITÄT BONN

# Motivation

# Motivation

- Photon Identification efficiency is important for: **inclusive prompt photon** and **di-photon** cross section, measurement of $H \rightarrow \gamma\gamma$ or any process that involves prompt photons in the final state.
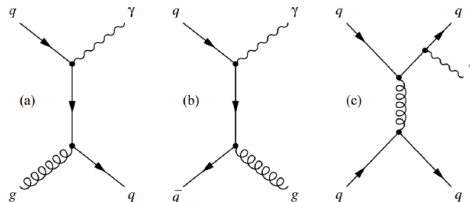


Figure 1: Prompt photon production via a)quark-gluon scattering b)quark-anti-quark annihilation c)bremsstrahlung radiation off of an outgoing quark.

- Fake photons consists of photons originated from the jets from the neutral hadron decays ($\pi^0$, $\eta$ mesons), and misreconstructed $e^- e^+$ pairs.

UNIVERSITÄT BONN

Figure 2: Canonical example of ABCD Method.

- **Matrix Method** is a data-driven method to measure the photon ID efficiency.

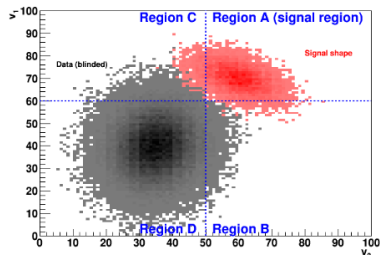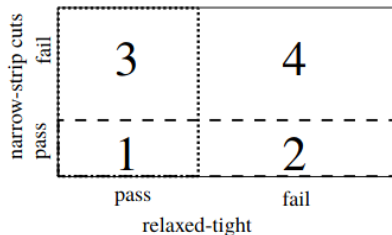  ▶ **Narrow-strip** and **relaxed-tight variables** as DVs (track isolation variables).

- Ideally, the matrix method is suitable for uncorrelated DVs. However, the photon isolation and shower shape variables are correlated.

- Instead of using rectangular cuts, **two neural networks** are used to determine the pass/fail statement for each narrow-strip and relaxed tight variables.

- Neural network outputs are used for the phase space separation in the matrix method.

UNIVERSITÄT BONN

# Matrix Method

# Matrix Method

- Matrix method uses narrow-strip, and relaxed tight variables.



- $\epsilon^{tight-ID} \equiv \frac{N_{ID}^s}{N^s}$

- $\hat{N}_a$: Number of track isolated photons in region a.

- $N_{ID}^T = N_{ID}^b + N_{ID}^s$

$$\epsilon^{tight-ID} = \frac{\frac{\epsilon_{ID}^{\hat{s}} - \epsilon_{ID}^{\hat{b}}}{\epsilon_{ID}^{\hat{s}} - \epsilon_{ID}^{b}} \cdot N_{ID}^T}{\frac{\hat{\epsilon} - \hat{\epsilon}^b}{\hat{\epsilon^s} - \hat{\epsilon^b}} \cdot N^T} \quad (1)$$

$$\hat{\varepsilon}_{ID}^b = \frac{\hat{N}_1^b}{N_1^b} \approx \frac{\hat{N}_3^b}{N_3^b} = \frac{R_p \cdot \hat{\varepsilon}_3 - A \cdot f_p \cdot \hat{\varepsilon}_3^s}{R_p - A \cdot f_p} \quad (2)$$
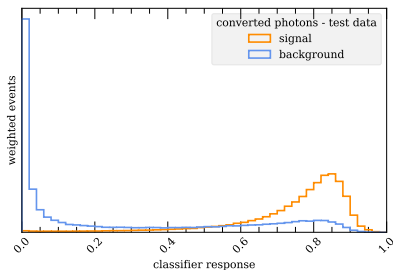
$$\hat{\varepsilon}^b = \frac{\hat{N}_{1+2+3+4}^b}{N_{1+2+3+4}^b} \approx \frac{\hat{N}_{2+3+4}^b}{N_{2+3+4}^b} = \frac{R_a \cdot \hat{\varepsilon}_{2+3+4} - A \cdot f_a \cdot \hat{\varepsilon}_{2+3+4}^s}{R_a - A \cdot f_a} \quad (3)$$

- This assumption leads to systematic uncertainties.(Correlation between narrow strip and track isolation.)
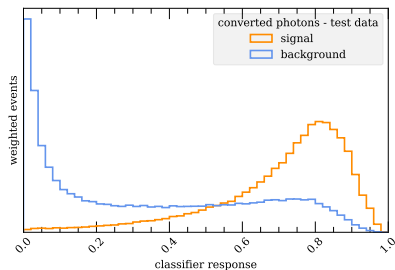
UNIVERSITÄT BONN

# Analysis

- In the analysis, the ATLAS data from 2015 to 2018, signal and background MC data are used.
- The data are put into the neural networks which outputs the pass/fail rate between 0-1.

Narrow Strip Variables

Relaxed Tight Variables

UNIVERSITÄT BONN

# NN Parameters

**Narrow Strip NN**

- Training Variables: $w_{\eta_1}$, $\Delta E_s$, $f_{side}$, $E_{ratio}$, $f_1$, $e_{277}$, $p_T$, $\eta$

  ECAL first layer

**Relaxed Tight NN**

- Training Variables: $R_\eta$, $R_\phi$, $w_{\eta_2}$, $s_{,tot}$, $R_{had}$, $R_{had_1}$, $p_T$, $\eta$

  ECAL second layer + Hadronic Leakage



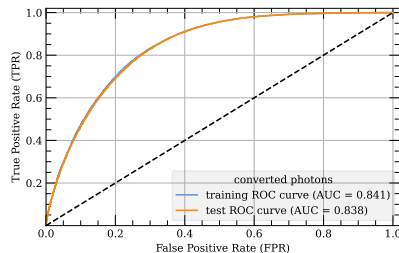Figure 3: ROC Curve of Narrow Strip NN



Figure 4: ROC Curve of Relax Tight NN

UNIVERSITÄT BONN

**The analysis script** classifies the events into 4-different regions according to the $y_{prediction}$.

- Converted photons
- Track Isolation Criteria:

  topoEtcone20 $< 6.5 \cdot 10^{-2} \cdot E_T$ and $\mathrm{ptCone}2\theta < 0.05 \cdot E_T$

- 4 $|\eta|$ intervals:

  $$[0.00, 0.60), [0.60, 1.37), (1.52, 1.81), [1.81, 2.37)$$

- 13 unequal $E_T$ bins from 25 to 1500 GeV.
- Pass threshold for Relax Tight: $y_{pred} > 0.55$
- Pass threshold for Narrow Strip: $y_{pred} > 0.6$

Calculate the following quantities,

- from the ATLAS data: $R_p, R_a, \hat{\epsilon}_{ID}, \hat{\epsilon}, \hat{\epsilon}_3, \hat{\epsilon}_{2+3+4}$
- from the Signal MC: $f_p, f_a, \hat{\epsilon}^s_{ID}, \hat{\epsilon}^s, \hat{\epsilon}^s_3, \hat{\epsilon}^s_{2+3+4}$

and save them into **json** files.

The script saves the root files into a Pandas dataframe, and creates dictionaries for each region. The data is stored for each $\eta$-$p_T$ bin.

UNIVERSITÄT BONN

# Systematic Uncertainties

- Sources of systematic uncertainties: MC statistics, track isolation requirement, detector material, closure uncertainty.

**Closure Test**

$$\Delta\hat{\varepsilon}^b_{ID} = \frac{\left|\hat{\varepsilon}^b_1 - \hat{\varepsilon}^b_3\right|}{\hat{\varepsilon}^b_1} \tag{4}$$

$$\Delta\hat{\varepsilon}^b = \frac{\left|\hat{\varepsilon}^b_{1+2+3+4} - \hat{\varepsilon}^b_{2+3+4}\right|}{\hat{\varepsilon}^b_{1+2+3+4}} \tag{5}$$

|  | converted photons | |
|---|---|---|
| $|\eta|$ interval | $\Delta\hat{\varepsilon}^b_{ID}[\%]$ | $\Delta\hat{\varepsilon}^b[\%]$ |
| $[0.0, 0.6)$ | 7.55 | 10.01 |
| $[0.6, 1.37)$ | 10.24 | 11.36 |
| $(1.52, 1.81)$ | 5.45 | 11.40 |
| $[1.81, 2.37)$ | 3.36 | 13.73 |

(6)

UNIVERSITÄT BONN

# Results

# Results for Track Isolation Efficiencies

To obtain the track isolation efficiency for the background sample:

$$a \cdot \left(\hat{\varepsilon}^b\right)^2 - b \cdot \hat{\varepsilon}^b + c = 0 \tag{7}$$

$$\begin{aligned}
a &= R_a - f_a \\
b &= R_a \cdot \left(\hat{\varepsilon}^S + \hat{\varepsilon}_{2+3+4}\right) - f_a \cdot \left(\hat{\varepsilon} + \hat{\varepsilon}^S_{2+3+4}\right) \\
c &= R_a \cdot \hat{\varepsilon}_{2+3+4} \cdot \hat{\varepsilon}^s - f_a \cdot \hat{\varepsilon}^S_{2+3+4} \cdot \hat{\varepsilon}
\end{aligned} \tag{8}$$

# Tight Identification Efficiency

$$\varepsilon^{tight-ID}\left(N_{ID}^T, N^T, \hat{\varepsilon}_{ID}^S, \hat{\varepsilon}^S, \hat{\varepsilon}_{ID}^b, \hat{\varepsilon}^b, \hat{\varepsilon}_{ID}, \hat{\varepsilon}\right) = \frac{\frac{\hat{\varepsilon}_{ID}-\hat{\varepsilon}_{ID}^b}{\hat{\varepsilon}_{ID}^s-\hat{\varepsilon}_{ID}^b} \cdot N_{ID}^T}{\frac{\hat{\varepsilon}-\hat{\varepsilon}^b}{\hat{\varepsilon}^s-\hat{\varepsilon}^b} \cdot N^T} \tag{9}$$

- High fluctuation due to **systematic errors**, and maybe due to a bug inside the code.

# Conclusion

# Conclusion

- The outputs of NNs are used instead of a cut based method.
- Matrix Method is an easy way to measure the photon ID efficiency.
- A further NN optimization should be done in the future for a better seperation, and for lower systematic uncertainties.
- A script that uses Uproot, Pandas, and numpy is written. Tensorflow is used for the neural network.
- Hands-on experience on working with real data, using modern ML techniques, a great motivation to start Master studies !...

That's it from my side, thanks for listening!

UNIVERSITÄT BONN

# Backup Slides

| Category | Variable name | Description |
|---|---|---|
| Hadronic leakage | $R_{\text{had}_1}$ | Ratio of $E_T$ in the first sampling layer of the hadronic calorimeter to $E_T$ of the EM cluster (used over the range $|\eta| < 0.8$ or $|\eta| > 1.52$) |
| | $R_{\text{had}}$ | Ratio of $E_T$ in the hadronic calorimeter to $E_T$ of the EM cluster (used over the range $0.8 < |\eta| < 1.37$) |
| ECAL first layer | $w_{\eta_1}$ | Lateral shower width, $\sqrt{\sum E_i (i - i_{\max})^2 / \sum E_i}$, where $i$ runs over all strips in a window of $3 \times 2$ $\eta \times \phi$ strips, and $i_{\max}$ is the index of the highest-energy strip calculated from three strips around the strip with maximum energy deposit |
| | $w_{s,\text{tot}}$ | Total lateral shower width $\sqrt{\sum E_i (i - i_{\max})^2 / \sum E_i}$, where $i$ runs over all strips in a window of $20 \times 2$ $\eta \times \phi$ strips, and $i_{\max}$ is the index of the highest-energy strip measured in the strip layer |
| | $f_{\text{side}}$ | Energy outside the core of the three central strips but within seven strips divided by energy within the three central strips |
| | $\Delta E_s$ | Difference between the energy associated with the second maximum in the strip layer and the energy reconstructed in the strip with the minimum value found between the first and second maximum |
| | $E_{\text{ratio}}$ | Ratio of the energy difference between the maximum energy deposit and the energy deposit in the secondary maximum in the cluster to the sum of these energies |
| | $f_1$ | Ratio of the energy in the first layer to the to the total energy of the EM cluster |
| ECAL second layer | $R_\eta$ | Ratio of the energy in $3 \times 7$ $\eta \times \phi$ cells over the energy in $7 \times 7$ cells centered around the photon cluster position |
| | $w_{\eta_2}$ | Lateral shower width, $\sqrt{\sum E_i \eta_i^2 / \sum E_i - (\sum E_i \eta_i / \sum E_i)^2}$, where $E_i$ is the energy and $\eta_i$ is the pseudorapidity of cell $i$ and the sum is calculated within a window of $3 \times 5$ cells |
| | $R_\phi$ | Ratio of the energy in $3 \times 7$ $\eta \times \phi$ cells over the energy in $3 \times 7$ cells centered around the photon cluster position |

UNIVERSITÄT BONN

**Classifier Settings/Hyperparameters:**

- 4 Layers with size of 32.
- Batch Size: 8192
- Epochs: 250
- Learning Rate: 0.001
- Adam Optimizer

UNIVERSITÄT BONN