

Moment Pooling: Gaining Performance and Interpretability Through Physics Inspired Product Structures

Rikab Gambhir With Athis Osathapan and Jesse Thaler

Email me questions at rikab@mit.edu! Based on [**RG**, Osathapan, Thaler, 23XX.XXXX (WIP)]



Typical Machine Learning Setup



Pictured: An Energy Flow Network (EFN):

$$\mathcal{O}(\{p_1,\ldots,p_M\}) = F\left(\sum_{i=1}^{M} z_i \Phi(\hat{p}_i)\right)$$

Typical Machine Learning Setup



The Moment-EFN

 $\mathcal{O}(\mathcal{P}) = F\left(\langle \phi^a \rangle_{\mathcal{P}}\right)$

Generalize!

EFNs^{*} can be thought of as taking the (weighted) **mean** of a latent particle representation ϕ – Let's generalize to *any* **moment**!

$\mathcal{O}_k(\mathcal{P}) = F_k\left(\langle \phi^a \rangle_{\mathcal{P}}, \langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}}, ..., \langle \phi^{a_1} ... \phi^{a_k} \rangle_{\mathcal{P}}\right)$

Claim: This "Moment-EFN" gives more efficient representations!

*Most of what I say here today also applies to Particle Flow networks or any other Deep-Sets inspired architecture!



The Moment-EFN (Details)



Claim: This "Moment-EFN" gives more efficient representations!



NSF AI Institute for Artificial Intelligence & Fundamental Interactions





Moment-EFNs



Using summary statistics of energy distributions within events to improve performance per parameter



e.g. Jet Angularities

In the moment language, even integer^{*} β jet angularities $\Leftrightarrow k = \beta^{\text{th}}$ moments!

$$\lambda^{(\beta)}(\mathcal{P}) = \sum_{i} z_{i} \left(\eta_{i}^{\beta} + \phi_{i}^{\beta} \right)$$
$$= \left\langle \eta^{\beta} \right\rangle_{\mathcal{P}} + \left\langle \phi^{\beta} \right\rangle_{\mathcal{P}}$$

For the normal (k = 1) EFN, this would require learning nonlinear functions!

Test: Train three networks to regress $\lambda^{(2)}$ from 100k QCD jet samples, with a latent dimension *L*:

- Linear Network: ϕ , *F* are 1 layer, linear functions, *L* = 2
- Small Network: ϕ , *F* are 2 layers, each with 4 nodes and *LeakyReLU*, *L* = 2
- "Large" Network: ϕ , F are 3 layers, each with 32 nodes and LeakyReLU, L = 8

Expect k = 2 to outperform k = 1 for smaller networks!

*Ask later about non-even or non-integer β angularities

e.g. Jet Angularities

In the moment language, even integer^{*} β jet angularities $\Leftrightarrow k = \beta^{\text{th}}$ moments!



e.g. Jet Angulari



Learns the simplest latent representations!

Training times are identical for k = 1 and 2!



10

Linearization

In principle, with a large enough k, we can approximate any^{*} observable with a **linear** F – like a Taylor expansion on distributions!

$$\mathcal{O}(\mathcal{P}) \approx F_a \langle \phi^a \rangle_{\mathcal{P}} + F_{a_1 a_2} \langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}} + F_{a_1 a_2 a_3} \langle \phi^{a_1} \phi^{a_2} \phi^{a_3} \rangle_{\mathcal{P}} + \dots$$

Intuitively: The product structure of moments capture all the nonlinearities that *F* would have captured

Can we understand the behavior of complex observables, like a Q/G discriminant^{**}, with just a few powers of k?

^{*}For well-behaved functions ^{**}For classification tasks, we linearize the log likelihood and apply a sigmoid or softmax at the end



e.g. Quark/Gluon Discrimination



Dataset of 500k Quark/Gluon Jets

Series of networks where ϕ has 3 layers and *F* is linear, but the width of ϕ and the latent dimension is sampled over – **ensemble of networks**

Consider k = 1, 2, 3, and 4

As with the previous example: Moments help achieve the same performance for less parameters!

See Backup slides for training details and dataset details



e.g. Quark/Gluon Discrimination

Compare to an arbitrary F network ...

13



Saturates at AUC = 0.88, consistent with <u>1810.05165</u>!

See backup slides for training details and dataset details. If we have time – see backup for performance versus latent dimension! Same performance for *lower* latent dimensions!

Interpretation

The linear approximation is not good enough here, unlike jet angularities ...



Four moments are not enough to build a linear approximation

\Leftrightarrow

Q/G discrimination cannot be easily captured by low-order moments! No simple closed-form approximation

Comparable with the story of Energy Flow Polynomials (EFPs): Degree 4 EFPS achieve an AUC \sim 0.75, but need d > 7 to achieve an AUC above 0.8



Conclusion $\mathcal{O}_k(\mathcal{P}) = F_k\left(\langle \phi^a \rangle_{\mathcal{P}}, \langle \phi^{a_1} \phi^{a_2} \rangle_{\mathcal{P}}, ..., \langle \phi^{a_1} ... \phi^{a_k} \rangle_{\mathcal{P}}\right)$

The Moment-Pooling structure improves open Deep Sets-type networks – better performance for smaller networks!

Can characterize how "simple" or "interpretable" an observable is by how many moments it takes to capture it in a linear basis – angularities are interpretable, Q/G discriminants may not be!



Rikab Gambhir – APS April – 17 April 2023

15

Appendices



16



Angularities Study (Details)

Dataset:

- 14 TeV Z+jet[g, uds] events generated in Pythia 8.226
- Jets clustered using AK4 (Fastjet 3.3.0)
- Keep p_T between 500 GeV and 550 GeV, |y| < 1.7
- 100k Train, 2.5k Val, 2.5k Test
- Angularities normalized to unit mean and standard deviation
- Particle p_{T} normalized to one.

Training:

- Batch Size: 512
- Epochs: 100
- Optimizer: Adam with learning rate 0.001



Q/G Study (Details)

Same dataset as angularity study, but with 500k training samples

For each of k = 1 ... 4:

- 1. Choose random integers F_{size} and ϕ_{size} from 1 ... 128, and *L* from 1 ... L_{max} , where L_{max} = 128 for k = 1, 64 for k = 2, 32 for k = 3, and 16 for k = 4.
 - a. Choose such that the number of network parameters is uniform in log scale.
 - b. For the linear *F* study, set F_{size} =1.
- 2. Initialize N = 3 Moment-EFNs (with different seeds), where the F and φ networks have three layers of the above size and latent dimension L.
 a. For the linear F study, instead choose the F network to be a single linear layer.
- 3. Train all *N* Moment-EFNs, using BCE loss, with the same hyperparameters are the angularities study. Record their AUCs.
- 4. Record the mean and standard deviation of the *N* AUCs and plot a single point. Repeat for 25 total points.



Quark/Gluon Discrimination: Latents

