ジェットの深層学習の決定過程の 解析と物理応用 (ミンコフキー汎関数を中心にして)

KEK 野尻美保子

ジェットの深層学習

- Jet のsubstructure からHeavy 粒子の選別 今後の高輝度化、高エ ネルギー化→ high pT のHiggs 粒子 top 粒子の選別, mass drop tagger など摂動的にも良くわかった量から出発
- BDT (機械学習)から深層学習
- 様々なアルゴリズムの提案 高次の量の利用→低次の量(Jet image など) MC レベルではよりよい結果を示す。例えば top Higgs vs
 QCD ジェットなど



安心できる量と不安な量

IRC SAFETY

- IRC safe object: soft or collinear emission に対して安定な量, subjet
- Soft collinear に対して不安定な量 number of tracks, particles
 MC modeling に大きな差 (いつまで立っても収束しないPythia と Hewig 実験データなど)
- →Jet Image の利用: 両者の間の区別が不明確 系統誤差の評価
 に問題はないか。実データで補正はやや楽観的ではないか。
- 理論研究:ジェットの構造の中でresummation でうまく記述できる 部分を切りだす。

深層学習によるジェット分類の中で、 IRC SENSTIVE な部分がどのように働いているか。

- 典型的なジェットイメージで分類に
- Higgs ハードな部分が複数ある(mass drop) vs QCD: あ まりない
- W, H, Z 粒子束が独立する傾向 vs QCD 撒き散らし。カ ラーコヒーレンスは以前からみたい量。
- クオーク (どちらかと言うと) 芯がある vs グルーオン より広い
- 注意 underlying event とかmultiple scattering とかあるのでLHC の環境では難しい量. Soft の部分を切り捨てて、摂動計算と比較できる物理をやりたいという流れもある。
- ・ 今のトーク ソフトな粒子の分布に注目したい。



Trimming, soft drop, Iterated soft drop



粒子分布を定量化する

ソフトなトラックの数だけではなくて、ソフトな粒子の相互の近接度を「定量的」に評価した
 い →キャリブレーションやMC turning に便利なように



 $N_1/N_0 = 16/9 = 1.78$

背景にある数学 INTEGRAL GEOMETRY

Object の相対的な位置を定量化する理論

凸体の測度(大きさ)の満たすべき性質

並進 M(gB) = M(B)加法 $M(B_1 \cup B_2) = M(B_1) + M(B_2) - M(B_1 \cap B_2)$ 連続性 $M(K_i) \to M(K)$ as $K_1 \to K$ for $K, K_i \in \mathcal{R}_i$

凸体の測度(大きさに関する定理)n次元では、この性質を満たす関数はn+1 しかない
 (ハドビガーの定理) これをミンコフスキー汎関数(MF)という。2次元であればこれ
 は 1. Surface Area (A), 2. Length of the boundary(L) 3. Euler characteristic (X)

点-> 点を中心に半径R の円を追加→ A(R), L(R),X(R) が点分布を記述

MFの他の物理応用

統計物理 左 多孔質体 真ん中: 微乳濁液 左 コロイド 体積の占有状況V, 表面の大きさ(S) 等に依 存して物性が変わる 図は Mecke and Stoyan (2000)

天文: 星の分布の定量化、銀河分
 布、シミュレーション結果の定量
 化、non-Gaussinaity of CMB,
 weak lensing..

点の集まりの意味を定量的に表現する時に使う



Kratochvil 1109.6334 Proving Cosmology with Weak Lensing Minkowski Functinal s



FIG. 1: Top left panel: example of a simulated 12-square-degree convergence map in the fulcial cosmology, with intrustic ellipticity noise from source galaxies and $\theta_G = 1$ arrain Gaussian smoothing. A source galaxy density of $\eta_{mal} = 15/4 \operatorname{rem}^{11}$ at redshift $z_* = 2$ was assumed. Other three panels: the excursion sets above three different convergence thresholds η_{max} , i.e. all packs with values above (below) the threshold are balax back (white). The threshold values are n = 0.0 (for pixel), n = 0.02 (bottom left), and n = 0.01 (bottom right). The Minkowski Functionals V_0 , V_1 , and V_2 measure the area, boundary length, and Euler characteristic (or genus), respectively, of the black regions as a function of threshold.

CNN はMF(ミンコフスキー汎関数) をみているか

- MF の一意性: 原理的には、ドットイメージとMF はある意味で等価。(情報は落ちてない)
- MF はCNNとは相性は良さそう。
- 基準点がないのでジェットアルゴリズムと相性が良い。そもそも、 Jet Algorithm の Voronoi 領域と"The Catchment Area of Jets "(Cacciari & Salam 2008)も同じ系列の数学
- 深層学習の分類問題との関係: QCD と t, Z, W の ジェットイメージのMF が十分に違っていれば、CNN はMF を学習しているだろう。



解析:TOP TAGGING の中で

Chakraborty, Lim, Nojiri, Takeuchi 2003.11787

- CNN (ベースモデル)を 高次量をインプットとするMLP と比較
- インプットの分割
 - S2: C correlator(Energy correlator) f(θ) =Ei Ej δ(θ-θij) 2 点関数 Tkachov (hep-ph 960138) Lim, Nojiri 1807.03312, Chakrabory, Lim Nojiri 1904.02092 ~ e₂^βの任意のβの情報を担っている。
 - Top なので、3点がメイン; : Leading subjetの粒子との2点関数
 +Leading subject の粒子を除いた2点関数
 - 上記を hard なsubjet 内の粒子に制限したもの S2trim (groom でも良い。
 - New :IRC sensitive なインプット (N0, N1, N0(pt> 4GeV), N1(pt> 4GeV)
 - 最終的なMLP は RN (relation network) + global な量(jet pt, mass, trimmed jet 情報と MF で作る





ROCはCNN と等価

- IRC insensitive な量だけ使った青い線と
 CNN は非常に差が大きい。
- Top jet vs QCD問題では、hit のあるカロリメータの数(N0)、その周りのエリア数を加えるとCNNと完全に同じになる
- Top ジェットの場合、カラーを持つ粒子のせいか、機械学習は ΔR= 0.1より遠距離のソフトな相関を学んでいない。
- Top がカラーを持っている。QCD 生成 など





Z boson efficiency

PREDICTABILITY

HiggsのT3カットへの応答を見る

N3LO 計算 (Mondini et al 1904.0896)

gloomed jet mass with NNLO+ N3LL resum e+e -> hadrons (factorization を保証)

 $\frac{\min[E_i, E_j]}{E_i + E_j} > z_{cut}$ うまく高次効果を計算できる フェーズスペースに特化した量を 計算する。

Kardos et al 2002.00942

mMDT groomed heavy hemisphere mass -2.10 $y_{cut}=0.1$ 0.35 -2.15-2.20 $z_{\rm cut} = 0.1$ 0.3 $Q = 91.2 \,\mathrm{GeV}$ -2.25 $\alpha_{\rm s}(m_Z) = 0.118$ -2.300.25 -2.35 $\frac{\rho}{\tau_n} \frac{d\sigma_g}{d\rho}$ 0.2 -2.40 $\mu = m_{\mu}$ -2.45 0.02 0.05 0.10 0.200.15 $\tau_3^{\rm cut}$ [GeV] 0.1□ NLO+N²LL QCD との関係がつくのはいずれにしても 0.05 N²LO+N³LL 3 POINT くらいまでで、後は実データ 0.0ベースで議論することになる。 0.01 0.02 $0.05 \ 0.1$ 0.2 0.5

LOCAL MINIMUM 問題の改善

CNN のloss function の最小化で、「真のminimum にたどり着くことはあまりない。ROC は安定しているが、個々のイベントに対して、違うseedで使ったclassifier は違う結果を出す。RN + MF は input が少ないのでevent ごとの結果も遙かに安定 (900-> 85)



N-subjettiness の場合 1807.04769



DARK JET の場合

Lim, Nojiri in preparation

- Dark Jet $pp \rightarrow Z' \rightarrow qD qD \rightarrow dark Parton shower \rightarrow \rho diag \rightarrow qq$
- カラーシングレットなシャワー:粒子がたくさんあるが、いくつかの
 カラーシングレットなクラスターがある状態



CNNの学習結果

mp=20GeV, 300GeV<pT<400GeV CNN のイベント選択は、
 MF(n>3) でカットをかけに行っていた。



MC TURNING とかCALIBRATION とか

- Top jet vs QCD ジェットは event generator が違うと 結構違う結果になる。
- ー番差にきいているのがQCDジェットの粒子数と広がり。結構違う、しかも実験データとも違う。
- [MC を実データで補正] MF の値が同じになるように、
 イベントにウエイトをつけると、一致がよくなる。





教訓とやれそうなこと

- CNN などのジェットイメージを使った訓練はインフラの物理の違い
 を使って分類を強化している。
- pixel wise な情報じゃなくて、MFのような「まとめ指標」の方が キャリブレーションにも便利かもしれない。
- MF: カラー構造の違う粒子の性質を効率的に捉えているように見える。
- N-subjettiness などの従来の指標ともコンシステント

おまけ NN のシステム



process pp →tt vs pp→2j 500GeV<pT<600GeV 150GeV<mj<200GeV

case 1

modulation for two point correlation two point correlation + Kin→5 outputs correlation to/without leading jet → 5 outputs

→ROC

case2 + N0(number of active pixel) →ROC

case 3 + N0,N1→ROC

Adding N₁ fill the gap between CNN and our approach.