

Summer School on Symmetries, Fundamental Interactions and Cosmology September 2018, Fraueninsel/Chiemsee

Matteo Cacciari LPTHE Paris and Université Paris Diderot

lets

Lecture I - Jet algorithms Lecture 2 - Jet substructure

[Includes material from Gavin Salam and Grégory Soyez]





Jet substructure

At the end of a jet finding (i.e. clustering) procedure, a jet is a **collection of constituents** to which we assign a 4-momentum (related to the sum of the 4-momenta of the constituents)

What is the **arrangement** of the constituents inside the jet?



Jet substructure

First studied by Mike Seymour in the early '90s to distinguish W jets from QCD jets

Topic revived about 10 years ago in order to study boosted objects



Jet substructure



The past ten years have seen en explosion in jet substructure studies, i.e. **how radiation is arranged within jets**, and what it can tell us



Convolved Feature Layers

Jet substructure

Jet declustering

Jet shapes

(calculate a function from radiation distribution



 $W \rightarrow WZ$ event

Convolutions



Max-Pooling

Why boosted objects



Heavy particle X at **rest**

Easy to resolve jets and calculate invariant mass, but signal very likely swamped by background (eg H→bb v.tt →WbWb)

Boosted heavy particle X

Cross section very much reduced, but acceptance better and some backgrounds smaller/ reducible



Mass of a single jet

Summing 'signal' and 'background' (with appropriate cross sections) shows how much the background dominates



Background only

nd onlySignal + backgroundPractically identical

This means that one can't rely on the invariant mass only. An appropriate strategy must be found to reduce the background and enhance the signal





Tagging and Grooming

The substructure of a jet can be exploited to

tag a particular structure inside the jet, i.e. a massive particle

▶ First examples: Higgs (2-prong decay), top (3-prong decay)

remove background contamination from the jet or its components, while keeping the bulk of the perturbative radiation (often generically denoted as grooming)

▶ First examples: filtering, trimming, pruning

Nomenclature

Groomer

procedure that always returns an output jet

(i.e. it only subtracts uncorrelated 'UE/pileup' radiation from it. This is used to 'clean' the jets from radiation largely unrelated to the fragmentation of the particle of interest)

Tagger

Procedure that might not return an output jet

(i.e. it either tags a heavy particle originating the jet or returns zero. This is used to identify a specific particle originating the jet.)

In practice, this classification is not always followed. In some cases it also denoted a 'tagger' a procedure that rejects background jets more often than signal jets

Why substructure

Scales: $m \sim 100 \text{ GeV}$, $p_t \sim 500 \text{ GeV}$

(e.g. electroweak particle from decay of ~ ITeV BSM particle)



need small R (< 2m/pt ~ 0.4) to resolve two prongs
need large R (>~ 3m/pt ~ 0.6) to cluster into a single jet

Possible strategies

- Use large R, get a single jet : background large
- Use small R, resolve the jets : what is the right scale?
 Also: small jets lead to huge combinatorial issues

Let an algorithm find the 'right' substructure

What jets to use for substructure?

Different jet algorithms will give different 'pictures' of what's inside a jet

Dendrogram

Used to represent graphically the sequence of clustering steps in a sequential recombination algorithm



Order of clustering here is A, B, C, D

The clustering sequence is 4-5 (A), 2-3 (B), 23-45 (C), 1-2345 (D)

First try

anti-kt

Hierarchical substructure







How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).

Second try

kt

Hierarchical substructure






How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).





How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).

 k_t clusters soft "junk" early on in the clustering

Its last step is to merge two hard pieces. Easily undone to identify underlying kinematics



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).

 k_t clusters soft "junk" early on in the clustering

Its last step is to merge two hard pieces. Easily undone to identify underlying kinematics



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).

 k_t clusters soft "junk" early on in the clustering

Its last step is to merge two hard pieces. Easily undone to identify underlying kinematics



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

This is crucial for identifying the kinematic variables of the partons in the jet (e.g. z).

 k_t clusters soft "junk" early on in the clustering

Its last step is to merge two hard pieces. Easily undone to identify underlying kinematics

This meant it was the first algorithm to be used for jet substructure.

Seymour '93 Butterworth, Cox & Forshaw '02

Third try

Cambridge/Aachen

Hierarchical substructure

Cambridge/Aachen



C/A distance measure

$$d_{ij} = \frac{\Delta y^2 + \Delta \phi^2}{R^2}$$

Cluster by merging the **closest** particles



















How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

C/A identifies two hard blobs with limited soft contamination



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

C/A identifies two hard blobs with limited soft contamination, joins them



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

C/A identifies two hard blobs with limited soft contamination, joins them



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?



How well can an algorithm identify the "blobs" of energy inside a jet that come from different partons?

C/A identifies two hard blobs with limited soft contamination, joins them, and then adds in remaining soft junk

The interesting substructure is buried inside the clustering sequence — it's less contamined by soft junk, but needs to be pulled out with special techniques

Butterworth, Davison, Rubin & GPS '08 Kaplan, Schwartz, Reherman & Tweedie '08 Butterworth, Ellis, Rubin & GPS '09 Ellis, Vermilion & Walsh '09

Hierarchical substructure



Slide by Gavin Salam
The IRC safe algorithms

	Speed	Regularity	UE contamination	Backreaction	Hierarchical substructure
kt	0000	\mathbf{T}	ŢŢ		00
Cambridge /Aachen	0000	Ţ	Ţ		000
anti-k _t	0000	00	♣/☺	00	×
SISCone	\odot		00		×

Array of tools with different characteristics. Pick the right one for the job

Matteo Cacciari - LPTHE

PRISMA Summer School - September 2018

QCD v. heavy decay

A possible approach for reducing the QCD background is to identify the two prongs of the heavy particle decay, and put a cut on their momentum fraction



Splittings and distances



For a given mass, the **background** will have smaller distance d_{ij} than the signal, i.e. it will tend to **cluster earlier** in the k_t algorithm

Potential tagger: last clustering in kt algorithm

This is where the hierarchy of the k_t algorithm becomes relevant. QCD radiation is clustered first, and only at the end the symmetric, large-angle splittings due to decays are reclustered

Alternative algorithms

- Suppose that for some reasons (which will become clearer later) one does not with to use the k_t algorithm
 - One must then find a way to determine what the **relevant splitting** (i.e. the one due to the decay, not to QCD radiation) is.

A possible approach is to use a Mass-Drop requirement: the clustering is **progressively undone**, and a splitting is the relevant one if both subjects are much less massive than their combination

$PP \rightarrow ZH \rightarrow v\bar{v}b\bar{b}$ The BDRS tagger/groomer

Butterworth, Davison, Rubin, Salam, 2008



A two-prong tagger/groomer for boosted Higgs, which

- Uses the **Cambridge/Aachen** algorithm (because it's 'physical')
- Employs a Mass-Drop condition, as well as an asymmetry cut to find the relevant splitting (i.e. 'tag' the heavy particle)
- Includes a post-processing step, using 'filtering' (introduced in the same paper) to clean as much as possible the resulting jets of UE contamination ('grooming')

BDRS: tagging

→ZH → vībb PP



BDRS: tagging

 $ZH \rightarrow vvbb$ PP



BDRS: tagging

 \rightarrow ZH \rightarrow vvbb PP



[NB. Parameters used $\mu = 0.67$ and $y_{cut} = 0.09$]

BDRS: filtering

 \rightarrow ZH \rightarrow vvbb PP



BDRS: filtering

$pp \rightarrow ZH \rightarrow vvbb$



BDRS: filtering

→ZH → vvbb PP



The low-momentum stuff surrounding the hard particles has been removed

Visualisation of BDRS

$pp \rightarrow ZH \rightarrow v\bar{v}b\bar{b}$

Butterworth, Davison, Rubin, Salam, 2008



Cluster with a large R

Undo the clustering into subjets, until a large asymmetry/mass drop is observed: tagging step Re-cluster with smaller R, and keep only 3 hardest jets: grooming step

First taggers/groomers

Mass Drop + Filtering

Butterworth, Davison, Rubin, Salam, 2008

Decluster with mass drop and asymmetry conditions Recluster constituents into subjets at distance scale R_{filt}, retain n_{filt} hardest subjets

Jet 'trimming'

Recluster constituents into subjets at distance scale R_{trim} , retain subjets with $p_{t,subjet} > \epsilon_{trim} p_{t,jet}$

Jet 'pruning'

S. Ellis, Vermilion, Walsh, 2009

Krohn, Thaler, Wang, 2009

While building up the jet, discard softer subjets when $\Delta R > R_{prune}$ and min(pt1,pt2) < ϵ_{prune} (pt1+pt2)

Aim: limit contamination from QCD background while retaining bulk of perturbative radiation

Trimming and pruner are a priori groomers, but can become taggers when combined with an invariant mass window test (if you can groom everything then there's no heavy particle in the jet)

The jet substructure maze



Soft Drop declustering

Larkoski, Marzani, Soyez, Thaler, 2014

Decluster and drop softer constituent unless Soft Drop Condition: $\frac{\min(p_{T1}, p_{T2})}{p_{T1} + p_{T2}} > z_{cut} \left(\frac{\Delta R_{12}}{R_0}\right)^{\beta}$

i.e. remove wide-angle soft radiation from a jet

The paper contains

- \checkmark analytical calculations and comparisons to Monte Carlos
- \checkmark study of effect of non-perturbative corrections
- \checkmark performance studies



Example of SoftDrop performance when used as a boosted W tagger

Alternatives to hierarchical substruct.

- If what we are interested in is the structure of the constituents of a jet, the "jet" itself is not the most important feature.
- A different algorithm, or simply the study of the constituents in a certain patch will also do. Selected alternatives are:
 - ▶ Use of jet-shapes to characterise certain features
 - e.g. *N-subjettiness*: how many subjets a jets appears to have

Thaler, van Tilburg, 2011

- Alternative ways of clustering
 - e.g. Qjets: the clustering history not deterministic, but controlled by random probabilities of merging. Can be combined with, e.g. pruning

Ellis, Hornig, Roy, Krohn, Schwartz, 2012

- ▶ Use information from matrix element
 - e.g. shower deconstruction: use analytic shower calculations to estimate probability that a certain configuration comes from signal or from background
- Use event shapes mimicking jet properties
 - e.g. JetsWithoutJets, mimicking trimming

Bertolini, Chen, Thaler, 2013

N-subjettiness

Thaler, van Tilburg, 2010

$$\tau_{N}^{(\beta)} = \sum_{i} p_{Ti} \min \left\{ R_{1,i}^{\beta}, R_{2,i}^{\beta}, \dots, R_{N,i}^{\beta} \right\}$$
Sum over constituents of a jet Distances to axes of N subjets

 T_N measures departure from N-parton energy flow: if a jet has N subjets, T_{N-1} should be much larger than T_N

N-subjettiness

Thaler, van Tilburg, 2010



 \mathcal{O}

Larkoski, Salam, Thaler 2013

Energy correlation functions Probes of N-prong structures without requiring identification of subjets

$$ECF(N,\beta) = \sum_{i_1 < i_2 < \dots < i_N \in J} \left(\prod_{a=1}^N p_{T_{i_a}} \right) \left(\prod_{b=1}^{N-1} \prod_{c=b+1}^N R_{i_b i_c} \right)^{\beta}$$

Angular (y-φ) distances between constituents

ECF(N+1) is zero if there are only N particles

More generally, if there are N subjets one expects ECF(N+1) to be much smaller than ECF(N) [because radiation will be mainly soft/collinear to subjets]

Discriminators

$$r_N^{(\beta)} \equiv \frac{\mathrm{ECF}(N+1,\beta)}{\mathrm{ECF}(N,\beta)}$$

small for N prongs: if N hard partons, small if radiation only soft-collinear

$$C_N^{(\beta)} \equiv \frac{r_N^{(\beta)}}{r_{N-1}^{(\beta)}} = \frac{\text{ECF}(N+1,\beta) \text{ECF}(N-1,\beta)}{\text{ECF}(N,\beta)^2}$$

A jet with a **small** C_N is more likely to have N prongs and at most soft/coll radiation



Note different values of β (chosen to maximise discriminating power)

The D functions are variations of the C ones

 $C_2^{(\beta)} = \frac{e_3^{(\beta)}}{(e_3^{(\beta)})^2} \qquad C_3^{(\beta)} = \frac{e_4^{(\beta)}e_2^{(\beta)}}{(e_3^{(\beta)})^2}$ Instead of $D_{2}^{(\beta)} = \frac{e_{3}^{(\beta)}}{(e_{2}^{(\beta)})^{3}} \qquad D_{3}^{(\alpha,\beta,\gamma)} = \frac{e_{4}^{(\gamma)} \left(e_{2}^{(\alpha)}\right)^{\frac{\gamma}{\alpha}}}{\left(e_{2}^{(\beta)}\right)^{\frac{3\gamma}{\beta}}} + x \frac{e_{4}^{(\gamma)} \left(e_{2}^{(\alpha)}\right)^{\frac{\gamma}{\beta}-1}}{\left(e_{2}^{(\beta)}\right)^{\frac{2\gamma}{\beta}}} + y \frac{e_{4}^{(\gamma)} \left(e_{2}^{(\alpha)}\right)^{\frac{2\gamma}{\alpha}-\frac{1}{\alpha}}}{\left(e_{3}^{(\beta)}\right)^{2}}$ define Top vs. QCD (Pythia 8) Attempt to improve the 160<mJ<240 GeV, pT>500 GeV, R=1.0 1.5 discriminating power, $D_{2}^{(\alpha,\beta,\gamma)}$ Relative Probability $(\alpha, \beta, \gamma) = (2, 0.8, 0.6)$ and to account for different 1.0OCD Jets Top Jets regions of phase space of radiation 0.5 [also, gives an idea of increasing 'sophistication', or complexification] 0.0

Matteo Cacciari - LPTHE

PRISMA Summer School - September 2018

2

3

 D_3

7

Robustness of substructure tools

Dasgupta, Fregoso, Marzani, Salam, 2013



Tools that are considered (or can be seen in Monte Carlo tests) to behave 'similarly' could cease to do so in different parameter regions

Analytic calculations of jet substructure

Dasgupta, Fregoso, Marzani, Salam, 2013



- Analytical understanding of 'kinks' in distributions
- Check of Monte Carlo predictions
- Other analytical investigations: Rubin 2010 (filtering), Walsh, Zuberi 2011 (jet substructure with SCET), Feige Schwartz, Stewart, Thaler 2012 (Nsubjettiness), Dasgupta, Marzani, Powling 2013 (groomed jet mass), ...

$$\frac{1}{\sigma} \frac{d\sigma}{dm^2}^{(\text{trim, LO})} = \frac{\alpha_s C_F}{\pi} \int_0^1 dz \, p_{gq}(z) \int \frac{d\theta^2}{\theta^2} \, \delta\left(m^2 - z(1-z)p_t^2\theta^2\right) \times \\ \times \left[\Theta\left(z - z_{\text{cut}}\right)\Theta\left(1 - z - z_{\text{cut}}\right)\Theta(\theta^2 - R_{\text{sub}}^2) + \Theta(R_{\text{sub}}^2 - \theta^2)\right] \Theta\left(R^2 - \theta^2\right)$$

Recap of Lecture 2

The big news of the past few years has been the development of taggers and groomers using properties of jet substructure, through

- declustering
- ▶ jet shapes
- direct analysis of images (machine learning)

These techniques have been commissioned by experimental collaborations proven their worth in 'Standard Model' analyses. They are now being implemented in BSM searches

Bonus: Jet areas

A jet's area is **defined** as the extent of the region where infinitesimally soft particles get clustered into the jet

More in details, a jet's **active area** is the extent of the region where a distribution of infinitesimally soft particles, that can also cluster among themselves, is clustered into the jet

A jet's active area, represented by the coloured regions, measures a jet's **susceptibility to contamination** from soft particles like underlying event and pileup



From jet 'reach' to jet areas

MC, Salam, Soyez, 0802. I 188

Not one, but three **<u>definitions</u>** of a jet's size:

Passive area

Place a single very soft particle (a '**ghost**') in the event, measure the extent of the region where it gets clustered within a given jet

Reach of jet for **pointlike** radiation

Active area

Fill the events with many very soft particles ('**ghosts**'), cluster them together with the hard ones, see how many get clustered within a given jet

Reach of jet for **diffuse** radiation

• Voronoi area

Sum of areas of intersections of Voronoi cells of jet constituents with circle of radius R centred on each constituent

Coincides with passive area for k_t algorithm

(In the large number of particles limit all areas converge to the same value)

Jet passive area

Passive Area

Add a **single** ghost^{*} particle to the event. Move it around. Check if it gets clustered in a given jet J.

$$a(J) \equiv \int dy d\phi f(g(y,\phi),J) \qquad \qquad f(g,J) = \begin{cases} 1 & g \in J \\ 0 & g \notin J \end{cases}$$

* ghost particle: particle with infinitesimally small momentum with respect to all other particles in the event (in practice, $O(10^{-100} \text{ GeV})$)

Matteo Cacciari - LPTHE

Jet active area

Active Area

Add **many** ghost particles in random configurations to the event. Cluster many times. Allow ghosts to cluster among themselves too. Count how many ghosts <u>on average</u> get clustered into a given jet J.



$$A(J) = \lim_{v_g \to \infty} \langle A(J | \{g_i\}) \rangle_g$$

Active area of jet J

Jet areas: the single hard particle case

It is worth noting that, for a jet made of a single hard particle, while passive areas are indeed πR^2 , active areas are not

Active areas	kt	Cam/Aa	SISCone	anti-k _t
< A >/πR ²	0.81	0.81	I/4	

Only anti-k_t has the behaviour one would naively expect, i.e. area = πR^2

Jet areas: the pure ghost jet case

Jets composed exclusively of ghosts can also be clustered, and their active area measured

Active areas	kt	Cam/Aa
< A >/πR ²	0,554	0,551

Active area distributions



For a roughly uniformly soft background, anti-k_t gives many small jets and many large ones (you can't fill a plane with circles!)

Active area: analytical results?

Can these active area distributions be calculated analytically?

At least the averages?

At least for pure ghost jets?

A challenge

Can we predict analytically the active area distribution or, at least, its average, for the **self-clustering of a large number of 'particles'**?

Consider a simplified toy-model: I-dimensional clustering, and an "euclidean" recombination scheme: the recombination of two 'particles' is simply the mid-point of their coordinate.

N 'particles' are distributed randomly over a length L, such that $T = \frac{T}{N} = \frac{T$

$$a \equiv L/N \ll 2R \ll L$$

where R is the 'radius' parameter of the Cambridge-like clustering algorithm

Upon clustering, we get n_j jets, with number of constituents n_c. The average 'area' (over many events) of a jet will be $\langle A \rangle = a \langle n_C \rangle$



(In this toy model, or in a similarly simple one.What about in higher dimensions?)