

Voyages into Machine Learning

Jesse Thaler

(1)

Voyages Before
the Sun

July 29

-July 6, 2019

based on: 1708.02949 with Eric Metodiev & Ben Nachman

1712.07124 with Eric Metodiev & Patrick Komiske
1810.05165

1802.00008 with Eric Metodiev
1809.01140 + Patrick Komiske

I was initially trained as a BSM model builder,
but over time I've become interested in methods to
let the data "speak for itself." This includes
model independent searches as well as a suite of
generally-applicable jet substructure techniques.

Ultimately, our field is moving towards "anomaly
detection", where we search for deviations from
the Standard model without specifying what BSM
physics we are looking for. I don't know how to
make that well defined, though, but I'm trying.
(and why I'm on this boat)

(2)

A big change in my research philosophy came with my grad students: Eric Metodiev & Patrick Komiske.

They have more of a CS / data science background, and they really pushed me to think about how to define problems explicitly.

ML \neq black box

ML = well-specified optimization problem.

But question is: what to optimize?!

(and subject to what constraints?)

More generally, ML yields functions defined by data.

Supervised learning: functions defined by labelled training data.

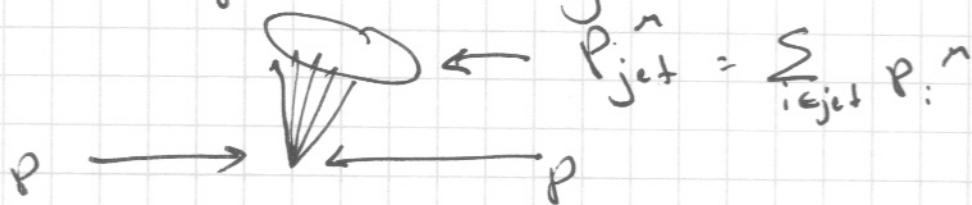
Unsupervised learning: functions defined by data at hand.

Many problems can be approached with this philosophy, but I'll focus on simplest: binary classification with case study of quark vs. gluon discrimination.

(3)

q vs. g : who cares?

First, quarks/gluons create jets in our detector



Showering, fragmentation, hadronization, detector effects.

$$q: C_F = 4/3$$



$$g: C_A = 3 \leftarrow \text{~twice as much radiation.}$$



If you can discriminate q vs. g (with some fidelity), you can ~~be~~ improve, e.g. SUSY searches
@ colliders

Signal: $pp \rightarrow \tilde{g} \tilde{g} \rightarrow q\bar{q} \rightarrow q\chi \bar{\chi}$ \Rightarrow 4 quark jets

Background: QCD multijet \Rightarrow many gluon jets.

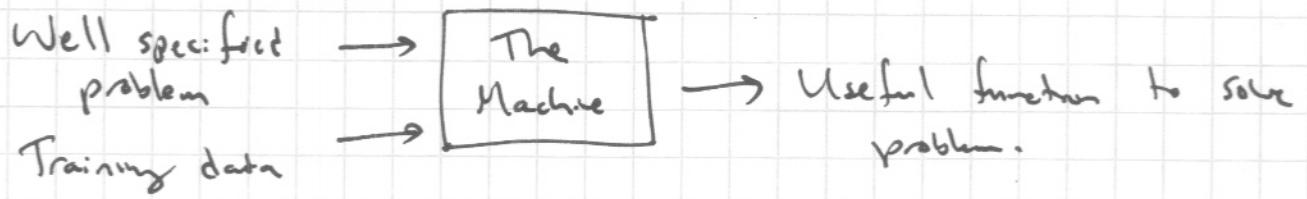
What is a quark or gluon jet?

color triplet vs. color octet, but
hadrons are color neutral (hmm...)

How do you optimally distinguish quark from gluon?
(assuming you could even define it...)

Remarkably, same answer: ML on mixed samples of
quarks & gluons. (with many caveats)

Cartoon of Machine Learning



For quark/gluon: we want $h(\text{jet}) = \begin{cases} 1 & \text{for quark} \\ 0 & \text{for gluon} \end{cases}$

We have unlabelled jet samples.

But we have strong theory knowledge / bias.

(5)

First, best you can ever do is

$$h_{\text{opt}}(\text{jet}) = \frac{p(\text{jet} | g)}{p(\text{jet} | g) + p(\text{jet} | b)}$$

(Neyman-Pearson lemma) Basically if there is a ~~big~~
chance to get same configuration from both quark/gluon,
you have irreducible overlap and can only guess.
(with some chance of being right)

But with enough labelled training data, you
can learn $h(\text{jet}) \approx h_{\text{opt}}(\text{jet})$ with gradient descent

$$\text{loss} = \underbrace{\left\langle (h(\text{jet}) - 1)^2 \right\rangle}_{\substack{\uparrow \\ \text{thing I'm} \\ \text{trying to} \\ \text{minimize.}}} + \underbrace{\left\langle (h(\text{jet}) - 0)^2 \right\rangle}_{\substack{\uparrow \\ \text{thing I'm trying} \\ \text{to determine.}}} \underbrace{1}_{\substack{\text{average over sample,} \\ \text{of pure jets.}}}$$

A few lines of algebra to show minimum
is $h_{\text{opt}}(\text{jet})$.

Modern ML = fancy functional choices for $h(\text{jet})$
that are easy to minimize.

(6)

What is $h(\text{jet})$?

"jet" = set of 4-vectors, with flow labels
 ↗
 unsorted, variable length list.
 ↗
 Quantum Mechanics!

A very powerful parametrization of $h(\text{jet})$ that respects permutation symmetry:

$$h(\text{jet}) : h\left(\underbrace{l_1, l_2, \dots, l_L}_{\text{"latent variables"}}\right)$$

$$l_i = \sum_j f(p_i^j)$$

↑ ↑
per particle

sum preserves
symmetry

Really simple! Possibly complete if L large enough.

Easy to plot f and h , can gain intuition for physics.

But need laptop for that. (show psychedelic images?)

(7)

Big question: how do I find $h_{\text{jet}}(\text{jet})$ if I don't have pure quark/gluon samples?

For me, this is when I began to appreciate that ML is not ^{really} about the machine.. It is about the physics.

Have mixed sample:

$$p(\text{jet}) = f_q p_q(\text{jet}) + f_g p_g(\text{jet})$$

$$\begin{array}{c} \uparrow \\ \text{Assume sample independence.} \end{array}$$

Assume you have access to samples with different quark/gluon ~~probabilities~~ fractions.

$$\text{Sample A: } f_q = 1 - f_g = f_A$$

$$\text{Sample B: } f_g = 1 - f_q = f_B$$

$$f_A p_q(\text{jet}) + (-f_A) p_g(\text{jet})$$

Optimally A vs. B : $\frac{(f_A + f_B) p_q(\text{jet}) + (2-f_A-f_B) p_{gum}(\text{jet})}{(f_A + f_B) p_q(\text{jet}) + (2-f_A-f_B) p_{gum}(\text{jet})}$

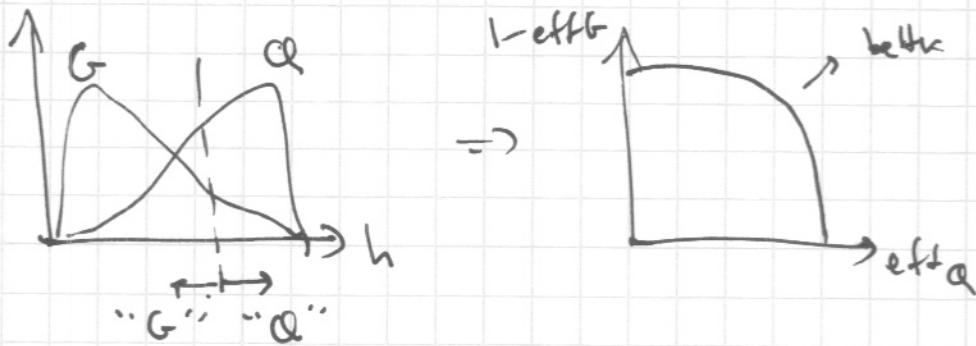
If you sit at this long enough...

(8)

h_{opt}^{mixed} is a function only of h_{opt}^{pure} ... and

$$\frac{\partial h_{opt}^{mixed}}{\partial h_{opt}^{pure}} = \frac{f_A - f_B}{(\text{Something})^2}$$

So h_{opt}^{mixed} is monotonically related to optimal classification!
(i.e. decision boundaries are the same)



Some ROC curve after monotone scaling.

So train on mixed samples \Rightarrow optimal decision boundaries
"CWoLa" = classification without labels.

We don't know purity of sample, though...

(i.e. what decision boundaries correspond to what purity)

But if we knew p_A and p_B we would be done!

7

Here is the magic: You can ~~use~~ classifier itself
 to define categories via "mutual irreducibility"

What is a "quark"? Hard to answer.

Is this a "quark"? Hard to answer.
 ↑
 particular jet

Is there any example jet (even if very rare), where
 I can say definitely it is a quark?
 (anchor bin)

Example from text processing:

"energy" physics or climate?

"energy conservation" physics or climate?

"Noether's Theorem",
 Kyoto Protocol



Identify "most quark-like" region

But we already have good discriminant! Cut
 really hard, and $\hat{h}_{opt} \rightarrow 1$ and $\hat{h}_{opt} \rightarrow 0$
 would be pure categories.

So basically endpoints or $\text{h}_{\text{opt}}^{\text{mixed}}$ can be solved to find f_{opt} (two equations, two unknowns)

Too good to be true? Yes, no starts at endpoints, but there are tricks to get around that.

What if you have three samples? Then you will find three "topics" (mutually irreducible distributions where each category is given in some "anchor bin")

If only two underlying categories, you'll find degenerate distributions (which you can quantify).

Bottom line: Given mixed samples, can ^{try} use ML
to define underlying categories.

Data-driven bootstrap. (aka. magic)

Looking ahead.. Anomaly Detection.

Want to find deviations from Standard Model.

If you have trustable MC, you can build

$$h_{\text{opt}}(\text{event}) : \frac{p(\text{event} | \text{data})}{p(\text{event} | \text{data}) + p(\text{event} | \text{MC})}$$

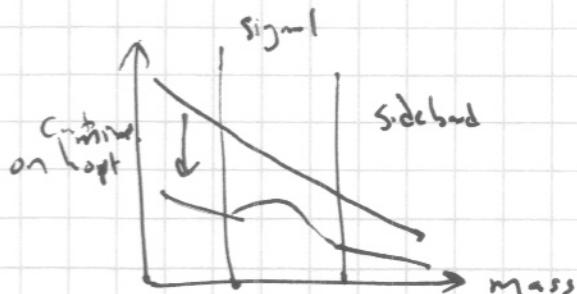
and any deviations from $1/2$ would be evidence for new physics

But we don't have a trustable MC always?

First attempt at solution:

CwLoa Hunting.

Collins, Hawke, Nachman 1805.02664



Build optimal signal / sideband
discriminant, cut hard
to maybe expose bump

This assumed we knew that new phys. would show up as bump. What if we don't know which feature separates signal from background?

(Hence 8 physicists on a boat)

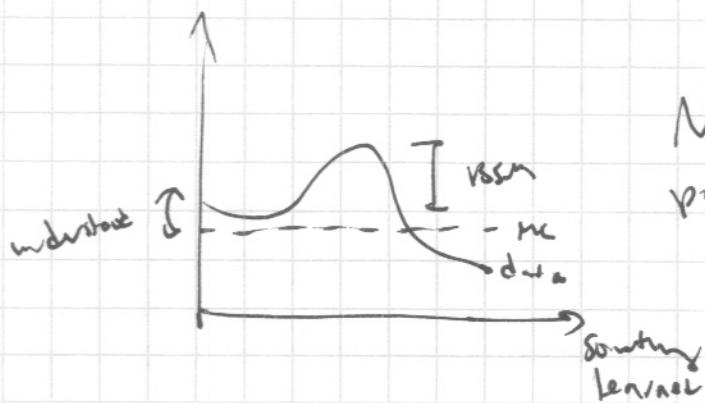
(12)

The problem: I have $p(\text{event} | \text{MC})$ but it isn't perfect.

I think new physics will show up in some "localized" ways as a general bump.

How do I find localized feature in $p(\text{event} | \text{MC}^{\text{data}})$ that differs from $p(\text{event} | \text{MC})$?

Localized, since less likely to be m.sunderstand MC effect.



Not so clear that this is possible, but now that I know top-down modeling works, inclined to try.