Bryan Ostdiek, TOP2020 Workshop September 17, 2020



Reducing uncertainty in top quark mass with Machine Learning



<u>Based on:</u> [<u>1804.00111</u>] T. Cohen, W. Hopkins, S. Majewski, and BO [<u>1909.09670</u>] T. Cohen, S. Majewski, BO, and P. Zheng [ongoing] F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz

bostdiek@g.harvard.edu

Splinters in the Stop Search





Splinters in the Stop Search







Czakon, Mitov, Papucci, Ruderman, and Weiler [1407.1043]



Czakon, Mitov, Papucci, Ruderman, and Weiler [1407.1043]

Now don't assume top mass; Cross section provides limit in stoptop plane



Czakon, Mitov, Papucci, Ruderman, and Weiler [1407.1043]

ATLAS used ratio of cross section at 7 TeV and 8 TeV as stop constraint (assuming top mass known)

 $pp \to \tilde{t}_1 \tilde{t}_1 \to e^{\pm} \mu^{\mp} \, b \, \bar{b} \, \tilde{\chi}_0^1 \, \tilde{\chi}_0^1 \, \nu \, \bar{\nu}$



T. Cohen, W. Hopkins, S. Majewski, and BO [1804.00111]

Machine learning the top quark mass



T. Cohen, W. Hopkins, S. Majewski, and BO [1804.00111]

Machine learning the top quark mass



T. Cohen, W. Hopkins, S. Majewski, and BO [1804.00111]

Machine learning the top quark mass







Eifert and Nachman [1410.7025]

- Previous constraints assumed top mass is known
- As stops can alter measured cross section, can also contaminate mass measurement



- Previous constraints assumed top mass is known
- As stops can alter measured cross section, can also contaminate mass measurement
- Different channels have varying amounts of contamination
- Reducing uncertainty can provide independent constraints
- Can machine learning help reduce uncertainty?

Template Method

1) Pick an observable which is correlated with the Monte Carlo top mass

2) Get distribution of observable at multiple values of the mass

3) Fit each distribution with a parametric function

4) Interpolate the function parameters as a function of the input mass

5) Go to real data and find the mass which fits the best

Template Method

1) Pick an observable which is correlated with the Monte Carlo top mass

0 *I:* Ratio of 3-jet invariant mass to 2-jet invariant mass
1 *I:* Common propagator value which maximizes combined likelihood
2 *I:* Invariant mass of the lepton and b-jet

2) Get distribution of observable at multiple values of the mass

- 3) Fit each distribution with a parametric function
- 4) Interpolate the function parameters as a function of the input mass
- 5) Go to real data and find the mass which fits the best

Template Method

Pick an observable which is correlated with the Monte Carlo top mass
 Get distribution of observable at multiple values of the mass



3) Fit each distribution with a parametric function

4) Interpolate the function parameters as a function of the input mass

5) Go to real data and find the mass which fits the best

Template Method

1) Pick an observable which is correlated with the Monte Carlo top mass

- 2) Get distribution of observable at multiple values of the mass
- 3) Fit each distribution with a parametric function



4) Interpolate the function parameters as a function of the input mass5) Go to real data and find the mass which fits the best

Template Method

1) Pick an observable which is correlated with the Monte Carlo top mass

- 2) Get distribution of observable at multiple values of the mass
- 3) Fit each distribution with a parametric function
- 4) Interpolate the function parameters as a function of the input mass



5) Go to real data and find the mass which fits the best

Template Method

- 1) Pick an observable which is correlated with the Monte Carlo top mass
- 2) Get distribution of observable at multiple values of the mass
- 3) Fit each distribution with a parametric function
- 4) Interpolate the function parameters as a function of the input mass
- 5) Go to real data and find the mass which fits the best



Comments

1) Unfortunately robust: the template/parametric fitting function does not need to be a good fit to the data



2) Without worrying about overall fit, become susceptible to contamination



Machine learning the top quark mass

Contamination Results



- The measured mass in each channel is different
- Channels with missing energy are affected more
- Largest effect comes when stops decay through off-shell top

Contamination Results



Stop contamination affects the three channels differently

Reducing uncertainty (and central values converging) will make the stop hypothesis inconsistent with the data

F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

• Template method finds parameterized model which best fits a distribution of a single observable

F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

- Template method finds parameterized model which best fits a distribution of a single observable
- Uncertainty due to ISR/FSR estimated by altering Monte Carlo parameters determining the change in extracted mass

F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

- Template method finds parameterized model which best fits a distribution of a single observable
- Uncertainty due to ISR/FSR estimated by altering Monte Carlo parameters determining the change in extracted mass
- Uncertainty can be reduced by jet calibration and removing extra radiation through grooming

F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

- Template method finds parameterized model which best fits a distribution of a single observable
- Uncertainty due to ISR/FSR estimated by altering Monte Carlo parameters determining the change in extracted mass
- Uncertainty can be reduced by jet calibration and removing extra radiation through grooming
- Use NN as universal function; makes it easy to use multiple observables and may be able to use information which grooming throws out

Machine learning the top quark mass

- Examine semi-leptonic $t\bar{t}$ production at $\sqrt{s}=13~{\rm TeV}$
- Use leptonic side to tag event, obtain mass measurement in hadronic side
- Anti-kt jets with R=0.5, $p_T^j > 30~{\rm GeV}$ $|\eta| < 2.4$
- 150 GeV < $m_{3j} < 200 \ {\rm GeV}$
- Use PYTHIA 8 with the A14 tunes:
 - 4 PDF sets
 - 5 families of tune variations
- Compute $\Delta m_t^{\rm MC} \equiv \Delta m_t^{\rm fit} \frac{m_t^{\rm MC}}{m_t^{\rm fit}}$, where $\Delta m_t^{\rm fit}$ is half of spread in fit mass using different tunes



F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Schwartz [in progress]

Machine learning the top quark mass



F. Flesher, K. Fraser, C. Hutchinson, BO, and M. Shwartz [in progress]

Machine learning the top quark mass

W cation and soft drop



W calibration and soft drop



W cation and soft drop



Machine learning the top quark mass

Deep neural networks using Classification for Tuning and Reweighting [Andreassen and Nachman [<u>1907.08209]]</u>

Use parameterized neural network classifier to learn likelihood ratio [use binary cross entropy loss function]

$$L(y, y^{p}) = \frac{-1}{n} \sum_{i=1}^{n} \left(y_{i} \log(y_{i}^{p}) + (1 - y_{i}) \log\left(1 - y_{i}^{p}\right) \right)$$

Classify as reference sample or any other mass/tune

Infer mass by computing the loss on unknown sample and reference sample

Use combination of high-level and low-level variables



Machine learning the top quark mass

Reference Set

 $m_t^{\rm MC} = 175 \,\,{\rm Gev}$ Color reconnection range = 1.71 $\alpha_{\rm s}^{\rm MPI} = 0.127$ Space shower $\alpha_s = 0.127$ Space shower $p_{t^0}^{\text{Ref}} = 1.51$ Space shower $p_{t,\text{damp}}^{\text{Fudge}} = 1.04$ Space shower $p_{t,\max}^{\text{Fudge}} = 0.88$ Time shower $\alpha_s = 0.124$ 10^6 events

Sampling Set

 $m_t^{\mathrm{MC}} \in (170 \text{ GeV}, 176 \text{ GeV})$ Color reconnection range $\in (1.67, 1.75)$ $\alpha_{s}^{\text{MPI}} \in (0.116, 0.136)$ Space shower $\alpha_s \in (0.1025, 0.1525)$ Space shower $p_{t^0}^{\text{Ref}} \in (1.415, 1.755)$ Space shower $p_{t,\text{damp}}^{\text{Fudge}} \in (0.715, 1.575)$ Space shower $p_{t,\max}^{\text{Fudge}} \in (0.75, 1.085)$ Time shower $\alpha_s \in (0.097, 0.153)$ 10^6 events

ParticleFlowNetwork on jet constituents helps network adapt to different tune parameters





Machine learning the top quark mass



Machine learning the top quark mass

Bryan Ostdiek



Mass and tune parameters extracted by gradient descent minimization



- Soft drop makes inference more challenging
- Similar to profiling over nuisance parameters
- Training is finicky (1/5 random initializations do not learn)
- Empirically find that only works for a small mass window



Training data: 100k events for each $m_t \in (165 \text{ GeV}, 180 \text{ GeV}) \text{ steps of } 0.2 \text{ GeV}$ tunes 19-32 Randomly select 30k events and train on the whole distribution

Test on $m_t = (172, 172.5, 173, 173.5, 174)$ GeV

Test set size: 400k (5 separate test sets per mass) Randomly select 30k -> get mass Repeat many times: mean is prediction







- Network is able to interpolate between masses
- Network learns distributions for the PDFS and tunes
- Grooming does not help
- Other observables did not help (2-jet mass useful for network to calibrate)

ОШН



- Network is able to interpolate between masses
- Network learns distributions for the PDFS and tunes
- Grooming does not help
- Other observables did not help (2-jet mass useful for network to calibrate)

ОШН



- Precise determination of top quark properties
 constrains BSM
- Machine learning can reduce uncertainty
 28-55% more than soft drop and W calibration
- Methods can be used for other measurements (cross section, width, etc)





- Precise determination of top quark properties
 constrains BSM
- Machine learning can reduce uncertainty
 28-55% more than soft drop and W calibration
- Methods can be used for other measurements (cross section, width, etc)





Loss options



The LogCosh loss has a shallower slope for small errors, allowing for easier minimization



Tuning with DCTR

Slide stolen from Anders Andreassen BOOST 2019 talk

• A well trained model $f(x, \theta)$ minimizes the loss for any θ

$$f(x,\theta) = \underset{f'}{\operatorname{argmax}} \sum_{i \in \theta_{\mathbf{0}}} \log f'(x_i,\theta) + \sum_{i \in \theta} \log \left(1 - f'(x_i,\theta)\right)$$

Tuning with DCTR

Slide stolen from Anders Andreassen BOOST 2019 talk

• A well trained model $f(x, \theta)$ minimizes the loss for any θ

$$f(x,\theta) = \underset{f'}{\operatorname{argmax}} \sum_{i \in \theta_{0}} \log f'(x_{i},\theta) + \sum_{i \in \theta} \log \left(1 - f'(x_{i},\theta)\right)$$

• Given a sample with θ_1 unknown (e.g. data), then

$$\theta^* = \underset{\theta'}{\operatorname{argmax}} \sum_{i \in \theta_0} \log f(x_i, \theta') + \sum_{i \in \theta_1} \log \left(1 - f(x_i, \theta')\right)$$

imples that $\theta^* = \theta_1$

Tuning with DCTR

Slide stolen from Anders Andreassen BOOST 2019 talk

• A well trained model $f(x, \theta)$ minimizes the loss for any θ

$$f(x,\theta) = \underset{f'}{\operatorname{argmax}} \sum_{i \in \theta_{0}} \log f'(x_{i},\theta) + \sum_{i \in \theta} \log \left(1 - f'(x_{i},\theta)\right)$$

• Given a sample with θ_1 unknown (e.g. data), then

$$\theta^* = \underset{\theta'}{\operatorname{argmax}} \sum_{i \in \theta_0} \log f(x_i, \theta') + \sum_{i \in \theta_1} \log \left(1 - f(x_i, \theta')\right)$$

imples that $\theta^* = \theta_1$

• Can find solution by gradient descending on θ'