Multivariate Techniques for $t\bar{t}H \rightarrow b\bar{b}$ Analysis in ATLAS

Ziyu GUO Supervisors: Yann COADOU, Thierry ARTIERES

Top LHC France @ LPSC Grenoble

April 25, 2019



1. $t\bar{t}H \rightarrow b\bar{b}$ analysis

2. MVAs in the analysis with 36.1 ${\rm fb}^{-1}$ in Run 2

3. Exploring deep learning in $t\bar{t}H \rightarrow b\bar{b}$ analysis Recurrent neural networks Parse tree Adversarial training Multiclassification

$t ar{t} H ightarrow b ar{b}$ analysis with 36.1 fb $^{-1}$ in Run 2 PhysRevD.97.072016

- Analysis challenges and strategy:
 - The systematic uncertainty is dominated by the uncertainty in the $t\bar{t} + b$ -jets modeling
 - region categorization including $t\bar{t} + b$ -jets enriched CR, simultaneous fit
 - Low $t\bar{t}H$ signal x-section, and large irreducible background of $t\bar{t}$ +jets
 - a BDT to separate signal from bkg, used as the discriminant for the profile likelihood fit in signal-enriched regions
 - Multiple jets and b-jets in the final state, difficult to be matched to partons correctly
 - using MVAs to reconstruct $t\bar{t}H$



- ► This study: 1/ channel, resolved events
- ▶ Pre-selection for training: (5jets, $\geq 4b$ -jets), (≥ 6 jets, $\geq 4b$ -jets)@85%

MVAs in the analysis with 36.1 fb^{-1} in Run 2

- Reconstruction step: find the correct association between jets and partons which originate from *H*/*top* decays
 - Reco. BDT: pick the combination with the highest BDT score as the correct matching among all possible combinations on t*t*H (see next slide)
 - Likelihood discriminant (LHD): build probability distribution function under the signal/background hypotheses using 1D variable distributions from all possible combinations
 - MEM: exploit the full matrix element calculation to separate the signal from the background.



Reconstruction BDT

- ► **Goal**: find the correct association between jets and partons which originate from H/top decays
 - Try all possible jet-parton matches, leading to multiple possible combinations
 - Train on $t\bar{t}H$ sample only
 - Signal: correct matches up to 1 jet from W mis-matched. Bkg: all other improper combinations
- Train two different BDTs:
 - BDT: using vars uncorrelated to Higgs
 - ▶ BDT withH: adding Higgs information, e.g. ∆R(b, b) from Higgs candidate
- Some input variables in $\geq 6j$ region:



Reconstruction BDT

- ► To reduce number of combinations, *b* quarks are only associated to *b*-jets.
- ▶ Previously, b-jets are tagged if passing a tighter fixed b-tagging efficiency threshold (WP: Working Point) → low stats, light quark often mis-tagged
- Using the pseudo-continuous (PC) b-tagging, and starting from loose WP, give more correct matchings
 - Sorting jets w.r.t. b-tagging weight, the leading 4 jets are considered as b-jets, and the leading 4 jets cannot be used as light jets.

Reco BDT with Higgs, PC vs. the fixed *b*-tagging:

- 42%, two *b*-jets from Higgs were truth matched
- 5% improvement with PC b-tagging
- Almost 10% improvement for all b-jets truth matching.

Figure: Reco BDT performance



Classification BDT

- ▶ 2 BDTs are trained for events having 5 and \geq 6 jets, both with \geq 4 *b*-jets @85% WP, across CRs and SRs
- ▶ Signal: tītH , Bkg: tīt
- Input variables:

Date /

- Global event kinematics
- PC b-tagging of jets
- Variables from reconstruction BDT: the vars built based on the combination with the highest Reco BDT score.
- Outputs of reco BDT, LHD, MEM (where available)





MEM discriminant



Classification BDT

- Inclusive training: in 4b@85%
- Dedicated training: in 4b@60%, a dedicated BDT is trained with MEM as additional input variable
- Using the PC b-tagging, the inclusive training gives similar performance to the dedicated training. Thus the inclusive training is used in all SRs except in 4b@60%
- Intermediate MVAs output variables are dominant discriminants

	AUC	improve (%)
Kinematics	0.738	-
Kins + RecoBDT	0.756	2.4%
Kins + LHD	0.763	3.4%
Kins + RecoBDT + LHD	0.768	4.1%

 Adding MEM gains a little bit, correlated with LHD

Correlation Matrix (signal) ATLAS Simulation work in progress																	
	Linear correlation coefficients in %																
MEN_D1	12	-10		-7	40	-18	12	-7			-11		31	48	100		100
LHD_Disoriminant	-7				52								33	100	48	-	80
THRees_withH_best_TTHRees_withH													100	33			60
TTHRee, with Lost, Hypether, dR												100					60
TTHReco_best_Higgsbleptop_mass	44										100					-	40
TTHReco_best_Higgsleptop_dR	5					9				100	32	21			-6		20
TTHReco_best_bbHggs_dR				42				46	100	5					2		20
TTHReco_best_Higgs_mass	38					12	-6	100	46						-7	-	0
Apianarity_jeta							100								12		00
dEnj_MaxEta	11					100									-18		-20
nHggbb030_5or4					100	-7								52	40		-40
dRd_avg_SoH	29			100	-4	31			42								~~
H.a			100	-10		-29									7		-60
dRbb_MarPt_Sot4	5	100	-6	39		18									-10		-80
Mbb_MincR_Sort4	100	5		29		11	-4	38			44				12		100
He was and a first a series and a series of the series of																	
A Company of the second s																	

RNN motivation

- Reconstruction step: use similar information but from different aspects
 - MEM: super computationally expensive to run, so only built in one signal rich region.
 - Reco. BDT: takes into account variable correlations. But only use the combination with the highest score, and the truth matching rate is limited.
 - Likelihood discriminant (LHD): use all combinations, but not the correlations between variables.
- Classification step: a BDT using info from these reco-level MVAs, classify ttH and tt
- Goal: take into account both variable correlations, and more possible combinations



RNN for $t\bar{t}H ightarrow bar{b}$

- Classify $t\overline{t}H$ and $t\overline{t}$ events
- Each event is a sequence, using combinations as frames
- Each frame is represented by the input features of the corresponding combination
- ► Same input features as reconstruction BDT + classification variables



This way, RNN takes into account more info than reco BDT and LHD: more combinations and proper feature correlations.

Training setup

- Machine learning tools:
 - ► Keras, Tensorflow, scikit-learn, etc: data science libraries
 - uproot: stream ROOT data into ML libraries
 - Trained on GPU
- Train, validation, test splitting



To avoid over-fitting, monitor the AUC during the training, pick the training epoch with the largest validation AUC.



・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

RNN performance

- ▶ Models optimized: input \rightarrow 1 RNN layer with 60 cells (dropout values=0.4) \rightarrow output layer with 1 neuron
- Same inputs as Reco BDT (with Higgs info) and classification BDT, w/o intermediate MVAs. Listed in backup.
- ▶ RNN (0.790) as good (even slightly better) AUC perf. as BDT (0.789)
- Statistical effect: not much to gain with larger $t\bar{t}$ sample, same for $t\bar{t}H$



 In 1 step, RNN achieves equivalent performance to 2-step MVAs consisting of various techniques, using the same info.

Parse tree: motivation

- ► Goal: exploring models trained with low level features
- Combining our domain knowledge + neural networks.
 - Inspired by QCD-Aware Recursive Neural Networks for Jet Physics
- Analogy between parse tree and Feynman diagram
 - Design a tree structure analogous to physical process
 - Use 4-vector and b-tagging of 8 objects as input: 6 jets + lepton and neutrino
 - ► Go through from the leaves to the collision node, embedding the input space to another n-dimensional space.



Parse tree model

- Use the tree embedding representation for each combination, making up the sequence input for RNN.
- ► Also add in high-level inputs: calculated features without combination info, used by BDT: ΔR_{bb}^{arg} , $\Delta R_{bb}^{max,pT}$, $\Delta \eta_{jj}^{max\Delta\eta}$, $m_{bb}^{min\Delta R}$, N_{30}^{Higgs} , H1, Aplanarity



Tree+RNN performance

- Previous BDT and simple RNN are trained with high level features (e.g. combination based and global kinematics)
 - Comparable AUC: BDT (0.789), simple RNN (0.790)
- ▶ RNN with (4-vectors, PC *b*-tagging), worse AUC (0.781)
- Tree+RNN could learn useful info from low level features, almost no gains from high level vars
 - ► Tree+RNN (0.788) with (4-vectors, PC *b*-tagging) as good as BDT/RNN
 - Tree+RNN (0.789) with (4-vectors, PC *b*-tagging, global kinematics), as good as BDT/RNN.
- Replace the tree embedding with a classical DNN. AUC:
 - 0.776 with (4-vectors, PC b-tagging)
 - 0.783 with (4-vectors, PC b-tagging, global kinematics)
- Tree performance is always better than DNN: tree structure helps to learn from low level features.



Adversarial training

- ▶ Dominant impact on µ is from tt
 + ≥ 1b shape difference of nominal and systematic samples.
- Difference exists in the nominal and syst samples, but small (but quite large compared to ttH presence):

	AUC AUC	Separation
	trained on nomina	l only
rnn nominal	0.787 ± 0.001	184.416 ± 0.426
rnn syst.	0.778 ± 0.001	222.155 ± 0.895
t	rained on nominal-	+syst.
rnn nominal	0.784 ± 0.001	184.965 ± 0.697
rnn syst.	0.778 ± 0.001	221.890 ± 2.102

 Goal: train a classifier insensitive to the difference between nominal and systematic samples. (Following paper: Learning to Pivot with Adversarial Networks)

bdt with nominal



rnn with nominal



tree+rnn with nominal



Adversarial training

- Idea: train a discriminator adversarially to restrict the classifier to have similar outputs for two different samples:
 - Minimax solution: $\hat{\theta}_{c}, \hat{\theta}_{d} = \arg\min_{\theta_{c}} \max_{\theta_{d}} \left(L_{c}(\theta_{c}) \lambda L_{d}(\theta_{c}, \theta_{d}) \right)$
 - L: loss function, measuring inconsistency between the prediction and label
 - θ : networks connection weights, $\lambda(>0)$: to be tuned



Alternative training:

- ▶ Train classifier θ_c , θ_d fixed: $t\bar{t}H$ vs $t\bar{t}$, nominal to be close to syst output
- Train discriminator θ_d , θ_c fixed: nominal vs syst output.
- Repeated till the θ_d unable to discriminate two samples.

Adversarial training

Experience we learn from lots of experiments with our analysis:

- ▶ The shape difference between samples is hard to be separated by discriminator.
- The performance metric we use before likelihood fits, AMS1¹sum
 - ► Take into account both signal presence and nominal-systematic difference.

 $\begin{aligned} \mathsf{AMS1} &= \sqrt{\sum_{i}^{N} 2((s_{i} + b_{i}) \ln(\frac{s_{i} + b_{i}}{b0_{i}}) - s_{i} - b_{i} + b0_{i}) + \frac{(b_{i} - b0_{i})^{2}}{\sigma_{bi}^{2}}}, \, \mathsf{N} = \mathsf{total} \\ \mathsf{number of bins} \\ \sigma_{bi} &= |b0_{i} - b1_{i}| \\ b0_{i} &= \frac{1}{2}(bi - \sigma_{bi}^{2} + \sqrt{(b_{i} - \sigma_{bi}^{2})^{2} + 4(s_{i} + b_{i})\sigma_{bi}^{2}}) \end{aligned}$

 Unstable value, even after rebinning to have more stats in each bin. Even sensitive to GPU randomness.



15

AUC vs training evolution





A robust RNN with adversarial networks?

- Not completely clear how much improvement this brings:
 - Output distributions show no clear evidence. Plots below are the RNN output w/ and w/o adversarial training.
 - AMS1 is improved, but with large uncertainty
 - BDT (trained on nominal only) AMS1: 0.752, AUC:0.789

AUC AMSI	AUC AMS1
$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$

w/ adversarial training



w/o adversarial training



Event categorization

- ► To improve the significance, events are categorized into orthogonal regions: signal regions (SR) and control regions (CR)
- Latest publication:
 - Splitting events into subsets w.r.t to jets pseudo-continuous b-tag score
 - Examining the bkg composition manually
 - Merge the subsets with similar bkg



イロト イポト イヨト イヨト

Multi-output classification

- Neural network classifier: a natural way to have multi-output
 - ▶ 1 RNN layer (50 cells) + Fully Connected (FC) layers + 4 output nodes (ttH, ttb, ttl, ttc)
- Model structure not optimized
- Balancing the ttH, ttb, ttc, ttl for training
- Could split further the ttb category



イロト イポト イヨト

Usage of multi-output

More studies are needed to find an optimal way to use the multi-output.

- Multi-dimensional cut
- Foam approach
- Following CMS strategy:
 - Categorize the event as ttH if the ttH node gives the highest score. Same for ttb, ttc, ttl
 - Using the output cells as region definition:
 - tth cell: signal rich region
 - ttb cell: ttb rich region
 - ► ...
 - Use output node associated to category for simultaneous fit

Background decomposition (ATLAS simulation work in progress)

Region definition based on multiclass output



tth

111 532.7

tb 530.0

tth 21.7

Default definition



Comparing two methods:

- The tth node gives similar number of $t\bar{t}H$ events, better s/b, s/ \sqrt{b} w.r.t to default inclusive SR
- SR1 and SR2 give less ttc & ttl constituents.

Less events in CR2 (ttc).

Summary

In the analysis with 36.2fb^{-1} of data, two-step MVAs are used

- Event reconstruction step and signal/bkg classification
- PC b-tagging variables are used in both steps

Exploring deep learning techniques for the full Run 2:

- Binary classification:
 - RNN: better explore reconstruction combinatorics due to recurrent structure, one-step classification
 - Parse tree: include physics domain knowledge while designing the neural networks, more efficient to learn from low level features
- Systematic uncertainty reduction:
 - ► Use an adversarial network that discriminates between tt models, to reduce modeling systematics. AUC pays the price, but AMS1 sum gains
- Region definition:
 - More flexible and natural with neural networks to have multi-outputs
- Complete fit studies with these new models are underway, preliminary results of simple RNN and multiclass are promising.

▲ロト ▲御 ト ▲ 臣 ト ▲ 臣 ト ○臣 - のへ()

Backup

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで

Reconstruction BDT inputs

Variable		on
Variable	≥ 6j	5j
Topological information from $t\bar{t}$		
t _{lep} mass	√	\checkmark
t _{had} mass	\checkmark	-
Incomplete thad mass	-	\checkmark
W _{had} mass	\checkmark	-
Mass of W_{had} and b from t_{lep}	\checkmark	\checkmark
Mass of W_{lep} and b from t_{had}	\checkmark	\checkmark
$\Delta R(W_{\text{had}}, b \text{ from } t_{\text{had}})$	\checkmark	\checkmark
$\Delta R(W_{\text{had}}, b \text{ from } t_{\text{lep}})$	\checkmark	\checkmark
$\Delta R(\text{lep}, b \text{ from } t_{\text{lep}})$	\checkmark	\checkmark
$\Delta R(\text{lep}, b \text{ from } t_{\text{had}})$	\checkmark	\checkmark
$\Delta R(b \text{ from } t_{\text{lep}}, b \text{ from } t_{\text{had}})$	\checkmark	\checkmark
$\Delta R(q_1 \text{ from } W_{\text{had}}, q_2 \text{ from } W_{\text{had}})$	\checkmark	-
$\Delta R(b \text{ from } t_{\text{had}}, q_1 \text{ from } W_{\text{had}})$	\checkmark	-
$\Delta R(b \text{ from } t_{\text{had}}, q_2 \text{ from } W_{\text{had}})$	\checkmark	-
min. $\Delta R(b \text{ from } t_{\text{had}}, q \text{ from } W_{\text{had}})$	\checkmark	-
min. $\Delta R(b \text{ from } t_{\text{had}}, q \text{ from } W_{\text{had}}) - \Delta R(\text{lep}, b \text{ from } t_{\text{lep}})$	\checkmark	\checkmark
Topological information from Higgs		
Higgs boson mass	√	\checkmark
Mass of Higgs and q_1 from W_{had}	\checkmark	\checkmark
$\Delta R(b_1 \text{ from Higgs}, b_2 \text{ from Higgs})$	\checkmark	\checkmark
$\Delta R(b_1 \text{ from Higgs, lep})$	\checkmark	\checkmark
$\Delta R(b_1 \text{ from Higgs}, b \text{ from } t_{\text{lep}})$	-	\checkmark
$\Delta R(b_1 \text{ from Higgs}, b \text{ from } t_{\text{had}})$	-	\checkmark

Input variables to the reconstruction BDT in the single lepton channel. The subscript had(lep) indicates the hadronically (leptonically) decaying W or t and q_i refers to quarks from W.

26 / 24

Э

classification BDT inputs

T Hysixev D.57.072010							
Variable	Definition	Reg	ion]			
variable	Denniuon	≥ 6j 5j		Variables from reconstruction BDT output			
General kin	ematic variables			BDT	BDT output	√*	√ *
ΔR_{bb}^{avg}	Average ΔR for all <i>b</i> -tagged jet pairs	√	~	m _H	Higgs boson mass	1	1
$\Lambda R^{\max p_T}$	ΔR between the two <i>b</i> -tagged jets with the	1	_	m _{H,bkp top}	Mass of Higgs boson and b-jet from leptonic top	1	-
Lan bb	largest vector sum p_T	•		$\Delta R_{\text{Higgs bb}}$	ΔR between b-jets from the Higgs boson	1	1
$\Delta \eta_{jj}^{\text{max} \Delta \eta}$	Maximum $\Delta \eta$ between any two jets	1	\checkmark	$\Delta R_{H,t\bar{t}}$	ΔR between Higgs boson and $t\bar{t}$ system	√*	√*
min AR	Mass of the combination of the two b-tagged			$\Delta R_{H,lep top}$	ΔR between Higgs boson and leptonic top	1	-
m _{bb}	jets with the smallest ΔR	v .	-	$\Delta R_{H,b_{had top}}$	ΔR between Higgs boson and <i>b</i> -jet from hadronic top	-	√*
$m_{ii}^{\min \Delta R}$	Mass of the combination of any two jets with the smallest ΔR	-	1	Variable fro	m Likelihood calculation		
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Number of h ist pairs with invesiont mass within			D	Likelihood discriminant	 ✓ 	1
N ₃₀ ^{Higgs}	30 GeV of the Higgs boson mass	1	\checkmark	Variable fro	m Matrix Method calculation		
Hhad	Scalar sum of jet p _T	-	1	MEM_{D1}	Matrix Method	~	-
1	A P between the lenter and the combination			Variables fr	om b-tagging		
$\Delta R_{lep-bb}^{min \Delta R}$	of the two <i>b</i> -tagged jets with the smallest ΔR	-	~	w_{μ}^{H}	Sum of binned b-tagging weights of jets	~	1
Aplanarity	1.5 λ_2 , where λ_2 is the second eigenvalue of the momentum tensor [91] built with all jets	1	1	\mathbf{B}_{j^3}	3 rd jet binned <i>b</i> -tagging weight (sorted by weight)	~	1
	Second Fox-Wolfram moment computed using			B _{/⁴}	4th jet binned b-tagging weight (sorted by weight)	V	1
H1	all jets and the lepton	~	~	B _{j5}	5th jet binned b-tagging weight (sorted by weight)	✓	1

Dhue Dev D 07 070016

Table 8: Input variables to the classification BDT in the single-lepton channel. For variables from the reconstruction BDT, those with a $^{+}$ are from the BDT using Higgs boson information, while those with no $^{+}$ are from the BDT without Higgs boson information. The *MEM_{D1}* variable is only used in the UPSR, while *b*-tagging weights are not used in this region (no information as they are all equal, by construction).

LSTM and multiclassification inputs for 6j case

Reconstruction info

	Topological information from $t\bar{t}$
	tlep mass
	Ihad mass
	Whad mass
	Mass of Whad and b from thep
15 RecoBDT	Mass of W_{lep} and b from t_{had}
inputs w/o	$\Delta R(W_{had}, b \text{ from } t_{had})$
Higgs info	$\Delta R(W_{had}, b \text{ from } t_{lep})$
	$\Delta R(\text{lep}, b \text{ from } t_{\text{lep}})$
	$\Delta R(\text{lep}, b \text{ from } t_{\text{had}})$
	$\Delta R(b \text{ from } t_{\text{lep}}, b \text{ from } t_{\text{had}})$
	$\Delta R(q_1 \text{ from } W_{had}, q_2 \text{ from } W_{had})$
	$\Delta R(b \text{ from } t_{had}, q_1 \text{ from } W_{had})$
	$\Delta R(b \text{ from } t_{had}, q_2 \text{ from } W_{had})$
	min. $\Delta R(b \text{ from } t_{had}, q \text{ from } W_{had})$
	min. $\Delta R(b \text{ from } t_{had}, q \text{ from } W_{had}) - \Delta R(\text{lep}, b \text{ from } t_{\text{lep}})$
RecoBDT	Higgs boson mass
innuts Higgs	Mass of Higgs and q_1 from W_{had}
info	$\Delta R(b_1 \text{ from Higgs}, b_2 \text{ from Higgs})$
	$\Delta R(b_1 \text{ from Higgs, lep})$
B-tagging	
weights	MVAreco_b1Higgs_pseudobtag_5bins
	MVAreco_b2Higgs_pseudobtag_5bins

Global classification BDT inputs

General kin	ematic variables				
ΔR_{bb}^{avg}	Average ΔR for all <i>b</i> -tagged jet pairs				
$\Delta R_{bb}^{\max p_T}$	ΔR between the two <i>b</i> -tagged jets with the largest vector sum p_T				
$\Delta \eta_{jj}^{\max \Delta \eta}$	Maximum $\Delta \eta$ between any two jets				
$m_{bb}^{\min \Delta R}$	Mass of the combination of the two <i>b</i> -tagged jets with the smallest ΔR				
N_{30}^{Higgs}	Number of <i>b</i> -jet pairs with invariant mass within 30 GeV of the Higgs boson mass				
Aplanarity H1	$1.5A_2$, where λ_2 is the second eigenvalue of the momentum tensor [91] built with all jets Second Fox–Wolfram moment computed using all jets and the lepton				

MVAreco_b2Higgs_pseudobtag_5bins MVAreco_blepTop_pseudobtag_5bins MVAreco_qhadW_pseudobtag_5bins_1 MVAreco_qhadW_pseudobtag_5bins_2 MVAreco_bhadTop_pseudobtag_5bins_2

> ◆□ ▶ ◆ 白 ▶ ◆ 王 ▶ ◆ 王 ▶ ◆ 王 ◆ ○ へ ^(*) 28 / 24

AMS1 and AUC stability

- Trials with 100 random seeds
- 100 times of experiments. train, val, test fixed, or train and val shuffled for each trial.
- ▶ AMS1 of BDT: 1.397.

Table: TransfoD AMS1 (mean/std(%))

	Fixed	Shuffled
train	.756±.146(19.3)	.840±.211(25.1)
val	$1.080 \pm .136(12.6)$.918±.187(20.4)
test	$1.765 \pm .304(17.2)$	$1.849 \pm .300(16.2)$

Table: Old transform AMS1 (mean/std(%))

	Fixed	Shuffled
train	.616±.109(17.7)	$.654 \pm .126(19.3)$
val	$1.035 \pm .146(14.1)$.840±.241(28.7)
test	$1.635 \pm .365(22.3)$	$1.670 \pm .433(25.9)$

Table: TransfoD: AMS2_diff (mean/std(%))

	Fixed	Shuffled
train	.115±0.016(13.9)	.108±.020(18.5)
val	.017±.016(94.1)	$.027 \pm .041(152)$
test	$.029 \pm .009(31.0)$.026±.009(34.6)

Table: Old transform: AMS2_diff (mean/std(%))

	Fixed	Shuffled
train	.084±0.015(17.9)	.080±.018(22.5)
val	.009±.013(144)	.015±.035(233)
test	.019±.006(31.6)	.019±.007(36.8)

Table: AUC0 (mean/std)

	Fixed	Shuffled
train	.814±.005	.813±.006
val	.793±.001	$.792 \pm .004$
test	.789±.001	$.788 \pm .001$

Table: AUC1 (mean/std)

	Fixed	Shuffled		
train	.792±.002	.792±.003		
val	.789±.002	.788±.007		
test	.778±.001	.778±.001		