Boosted decision trees: single top quark production evidence at DØ

Yann Coadou

CERN

Workshop on Top Physics: from the TeVatron to the LHC Grenoble, 20 October 2007



Wish I were here...





Poor Me

"I want to go to Grenoble to work on Saturday!"



Thank you Harrison!

Yann Coadou (CERN) — Boosted decision trees: DØ single top evidence

Top quark physics

- top quark discovered in 1995 by CDF and DØ at the Tevatron
- Heaviest of all fermions

- Couples strongly to Higgs boson
- So far only observed in pairs, only at the Tevatron



©B.Vachon

Single top quark production

• Never observed before: electroweak production



t-channel (tqb)



 $\begin{array}{l} {\sf Run \ II \ D} \mbox{\it 0:} < 4.4 \ pb \ (370 \ pb^{-1}) \\ {\sf Run \ II \ CDF:} < 3.2 \ pb \ (700 \ pb^{-1}) \end{array}$

(*) $m_t = 175 \text{ GeV}$, Phys.Rev. D70 (2004) 114012

Why do we care? — $|V_{tb}|$, new physics

- Has never been observed before!
- It should happen in SM
- First measurement of $|V_{tb}|$





Direct access to $|V_{tb}|$

$$V_{CKM} = \left(\begin{array}{ccc} V_{ud} & V_{us} & V_{ub} \\ V_{cd} & V_{cs} & V_{cb} \\ V_{td} & V_{ts} & V_{tb} \end{array}\right)$$

- In SM, from constraints on V_{td} and V_{ts} : $|V_{tb}| = 0.9991^{+0.000034}_{-0.000004}$
- New physics, e.g. 4^{th} generation: $0.07 < |V_{tb}| < 0.9993$

New physics

- s and t channel cross sections differently sensitive
- s-channel: charged resonances (heavy W' boson, charged Higgs boson, charged top pion, etc.)
- t-channel: new interactions (FCNC, 4th generation, etc.)

Why do we care? — Spin, Higgs, analysis techniques

Top quark spin

- Large mass ⇒ top quark decays before it can hadronize (no top jets)
- First chance to study a bare quark!

- Top polarization reflected in angular distributions of decay products
- SM predicts high degree of left-handed tops ⇒ possible sign of new physics, or help pin down what new physics

Higgs searches



- Important background to WH associated Higgs production
- As soon as we discover it, somebody will try to get rid of it....

Why do we care? — Spin, Higgs, analysis techniques

Top quark spin

- Large mass ⇒ top quark decays before it can hadronize (no top jets)
- First chance to study a bare quark!

- Top polarization reflected in angular distributions of decay products
- SM predicts high degree of left-handed tops ⇒ possible sign of new physics, or help pin down what new physics

Higgs searches



- Important background to WH associated Higgs production
- As soon as we discover it, somebody will try to get rid of it....

Advanced analysis techniques

- Test of techniques to extract small signal out of large background
- If tools don't work for single top, forget about the Higgs and other small signals
- If tools don't work at Tevatron, not much hope for LHC

Before we go on...



!!! VERY IMPORTANT !!!

Understand your inputs well before you start playing with multivariate techniques

Event selection



Event selection — Agreement before tagging

- Normalize W+jets and multijet to data before tagging
- Checked 90 variables, 4 jet multiplicities, electron + muon
- Good description of data



Agreement after tagging



- It took 8 months to get there...
- The signal is smaller than the background uncertainty
- Need to get smarter

Totals	2 Jets	3 Jets	4 Jets
Data	697	455	246
Total Background	685	460	253
Signal	36	20	6



Multivariate analysis techniques



- All three analyses have similar sensitivity and give compatible measurements
- Details about decision trees only today

Decision trees

- Machine-learning technique, widely used in social sciences
- Idea: recover events that fail criteria in cut-based analysis
- Start with all events = first node
 - sort all events by each variable
 - for each variable, find splitting value with best separation between two children (mostly signal in one, mostly background in the other)
 - select variable and splitting value with best separation, produce two branches with corresponding events ((F)ailed and (P)assed cut)
- Repeat recursively on each node
- Splitting stops: terminal node = leaf



• DT output = leaf purity, close to 1 (0) for signal (bkg)

Ref: Breiman et al, "Classification and Regression Trees", Wadsworth (1984)

Tree construction parameters

Normalization of signal and background before training

same total weight for signal and background events

Selection of splits

- Iist of questions (variable_i > cut_i?)
- goodness of split (next slide)

Decision to stop splitting (declare a node terminal)

- minimum leaf size (100 events)
- insufficient improvement from splitting

Assignment of terminal node to a class

- signal leaf if purity > 0.5
- background otherwise

Splitting a node

Impurity i(t)

- maximum for equal mix of signal and background
- symmetric in p_{signal} and P_{background}
- Decrease of impurity for split s of node t into children t_L and t_R (goodness of split): Δi(s, t) = i(t) - p_L · i(t_L) - p_R · i(t_R)
- Aim: find split s* such that:

$$\Delta i(s^*,t) = \max_{s \in \{\text{splits}\}} \Delta i(s,t)$$

- Maximizing Δ*i*(s, t) ≡ minimizing overall tree impurity
- Yann Coadou (CERN) Boosted decision trees: DØ single top evidence

- minimal for node with either signal only or background only
- strictly concave ⇒ reward purer nodes

Examples

$$\begin{array}{l} \text{Gini} = 1 - \sum_{i=s,b} p_i^2 = \frac{2sb}{(s+b)^2} \\ \text{entropy} = - \sum_{i=s,b} p_i \log p_i \end{array}$$



Decision Trees - 49 input variables

Object Kinematics

 $p_{T}(jet1)$ $p_{T}(jet2)$ $p_{T}(jet3)$ $p_{T}(jet4)$ $p_{T}(obtest1)$ $p_{T}(otbest1)$ $p_{T}(otbest2)$ $p_{T}(tag1)$ $p_{T}(untag1)$ $p_{T}(untag2)$

Angular Correlations

 ΔR (jet1, jet2) $\cos(best1, lepton)_{besttop}$ cos(best1,notbest1) cos(tag1,alljets)alljets $\cos(tag1, lepton)_{btaggedtop}$ cos(jet1,alljets)alljets $\cos(jet1, lepton)_{btaggedtop}$ cos(jet2,alljets)alljets $\cos(jet2, lepton)_{btaggedtop}$ $\cos(\text{lepton}, Q(\text{lepton}) \times z)_{\text{besttop}}$ cos(lepton_{besttop},besttop_{CMframe}) cos(lepton_{btaggedtop},btaggedtop_{CMframe}) cos(notbest,alljets)alliets cos(notbest,lepton) cos(untag1,alljets)alljets cos(untag1,lepton)

Event Kinematics

Aplanarity(alliets,W) M(W.best1) ("best" top mass) M(W,tag1) ("b-tagged" top mass) $H_{T}(\text{alljets})$ H_T (alljets-best1) $H_T(\text{alljets}-\text{tag1})$ $H_{T}(alliets, W)$ $H_{\tau}(\text{iet1.iet2})$ $H_T(jet1, jet2, W)$ M(alljets) M(alliets-best1) M(alliets-tag1) M(jet1,jet2) M(jet1, jet2, W) M_{T} (jet1, jet2) $M_T(W)$ Missing E_{T} p_T(alljets-best1) $p_T(alljets-tag1)$ p_T (jet1,jet2) $Q(lepton) \times \eta(untag1)$ $\sqrt{\hat{s}}$ Sphericity(alliets, W)

- Adding variables does not degrade performance
- Tested shorter lists, lost some sensitivity
- Same list used for all channels

Decision tree output

Measure and apply

- Take trained tree and run on independent pseudo-data sample, determine purities
- Apply to data
- Should see enhanced separation (signal right, background left)
- Could cut on output and measure, or use whole distribution to measure

Limitations

- Instability of tree structure
- Piecewise nature of output



Advantages

- DT has human readable structure (no black box)
- Training is fast
- Deals with discrete variables
- No need to transform inputs
- Resistant to irrelevant variables

Boosting a decision tree

Boosting

- Recent technique to improve performance of a weak classifier
- Recently used on decision trees by GLAST and MiniBooNE
- Basic principle on DT:
 - train a tree T_k
 - $T_{k+1} = \text{modify}(T_k)$

Ref: Freund and Schapire, "Experiments with a new boosting algorithm", in *Machine Learning: Proceedings of the Thirteenth International Conference*, pp 148-156 (1996)

AdaBoost algorithm

- Adaptive boosting
- Check which events are misclassified by *T_k*
- Derive tree weight α_k : $err_k = \frac{\sum_{i=1}^N w_i \times isMisclassified_k(i)}{\sum_{i=1}^N w_i}$ $\alpha_k = \beta \times \ln((1 - err_k)/err_k)$
- Increase weight of misclassified events by e^{α_k}
- Train again to build T_{k+1}
- Boosted result of event *i*: $T(i) = \sum_{k=1}^{N_{\text{tree}}} \alpha_k T_k(i)$
- Averaging \Rightarrow dilutes piecewise nature of DT
- Usually improves performance

Analysis validation

Ensemble testing

- Test the whole machinery with many sets of pseudo-data
- Like running DØ experiment 1000s of times
- Generated ensembles with different signal contents (no signal, SM, other cross sections, higher luminosity)

Ensemble generation

- Pool of weighted signal + background events
- Fluctuate relative and total yields in proportion to systematic errors, reproducing correlations
- Randomly sample from a Poisson distribution about the total yield to simulate statistical fluctuations
- Generate pseudo-data set, pass through full analysis chain (including systematic uncertainties)

All analyses achieved linear response to varying

input cross sections and negligible bias

Cross-check samples

- Validate methods on data in no-signal region
- "W+jets": =2jets, H_T(lepton, ∉_T, alljets) < 175 GeV
- "ttbar": =4jets, *H*_T(lepton,∉_T,alljets) > 300 GeV



Good agreement of model with data

Decision trees on data

- 36 different decision trees
 - 3 signals (s,t,s+t)
 - 2 leptons (*e*, *µ*)
 - 3 jet multiplicities (2,3,4 jets)
 - 2 *b*-tag multiplicities (1,2 tags)
- For each signal train against the sum of backgrounds



Results

 $\sigma_{s+t} = 4.9 \pm 1.4 \text{ pb}$ p-value = 0.035% (3.4 σ) SM compatibility: 11% (1.3 σ)

Evidence for single top production!





Comparison for DØ single top analysis



- Cannot know a priori which method will work best
- ullet \Rightarrow Need to experiment with different techniques

22

Conclusion

First evidence for single top quark production (D \emptyset decision trees)

 $\sigma(p\bar{p} \rightarrow tb + X, tqb + X) = 4.9 \pm 1.4 \text{ pb}$ 3.4 σ significance

First direct measurement of $|V_{tb}|$ (DØ decision trees)

 $\begin{aligned} |V_{tb}f_1^L| &= 1.3 \pm 0.2\\ \text{assuming } f_1^L &= 1: \quad 0.68 < |V_{tb}| \le 1 \text{ (0) 95\% CL}\\ \text{(Always assuming } V_{td}^2 + V_{ts}^2 \ll V_{tb}^2 \text{ and pure } V-A \text{ and CP-conserving } Wtb \text{ interaction}) \end{aligned}$

Published in Phys. Rev. Lett. 98, 181802 (2007) (hep-ex/0612052)

New preliminary combination of DT, ME and BNN

 $\sigma(p\bar{p} \rightarrow tb + X, tqb + X) = 4.7 \pm 1.3 \text{ pb}$

3.6 σ significance

• A lot more data already at hand

Single top prospects — Tevatron and LHC

Tevatron

- By 2009 we should have observed single top production and measured its cross section to 15-20%
- $|V_{tb}|$ is then known to ${\sim}10\%$

LHC

- Much larger production rates: $\sigma_s^{t/\bar{t}} = 6.6/4.1 \text{ pb } (\pm 10\%)$ $\sigma_t^{t/\bar{t}} = 156/91 \text{ pb } (\pm 5\%)$ $\sigma_{tW}^{t/\bar{t}} = 34/34 \text{ pb } (\pm 10\%)$ • Try to observe all three channels (s-channel challenging)
- $|V_{tb}|$ measured to percent level
- Large samples \Rightarrow study properties

