

# Estimation of btagging efficiency

Igor Marfin

August 17, 2010

Outline

Calibration  
methods

Likelihood ratio  
technique  
Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary

- 1 Calibration methods
  - Likelihood ratio technique
  - Flavor-tag consistency method
- 2 Kinematical Fit
- 3 Events preselection
- 4 Kinematical variables
- 5 Samples used for study
- 6 MVA Training
- 7 MVA Analysis
- 8 Estimation of btagging eff.
- 9 Summary



- The P<sub>tr</sub>el method. This method is based on measuring  $P_{tr}el = p_{\mu} \times p_{\mu+jet}/|p_{\mu+jet}|$  from events with two reco jets and one non-isolated muon before and after btagging. Then number of bjets before and after btagging can be fitted from P<sub>tr</sub>el distribution with MC templates.
- The System8 method. Based on the same events as before but taking into the account cut on P<sub>tr</sub>el and number of jets before and after btagging. Solving 8 equations on numbers of jets the performace is estimated.
- **Top-quark based method:** Likelihood technique. Using likelihood cut one can obtain semimuonic ttbar events with highly enriched b-jet content and suppressed background. Then the fraction of b-jets  $x_b = \frac{bjets}{alljets}$  is calculated before and after btagging. Using  $x_b$  and mistag rates (estiamted from MC) one can get btagging efficiency.
- **Top-quark based method:** Flavor-tag consistency method. The btag efficiency and mistag rates can be obtained from minimazing log-likelihood function  $L = 2 \log \prod_n P(N_n, \bar{N}_n)$ , where  $N_n$  and  $\bar{N}_n$  are the measured and expected number of events with  $n = 0, 1, 2$  tagged jets.  $P$  is the Poisson distribution.

This method is based on estimating the fraction of bjets in MC and measuring tagged jets in experiments after preselection and likelihood cut.

Likelihood function is used to suppress remaining background (mainly from  $W$ +jets events) and get events with the highly enriched b-jet content.

The **MVA technique** is used to construct the likelihood function.

- The preselection of events is done.
- The likelihood function  $L = \prod_i f_i(x_i)$  is constructed from MC, where  $x_i$  is some observable,  $f_i = \frac{S}{S+B}$  with  $S(B)$  -  $x_i$  distribution derived bin by bin way for Signal (Background).
- The fraction of bjets  $x_b = \frac{bjets}{alljets}$  is estimated from MC events survived the selection and the likelihood cut.
- The mistag rate  $\epsilon_0$  is estimated from MC.
- The fraction of tagged jets  $x_{tag} = \frac{tagjets}{alljets}$  is **measured from data** passed through the selection and the likelihood cut.
- One can calculate btagging efficiency as  $\epsilon_b = (x_{tag} - \epsilon_0 * (1 - x_b)) / x_b$ .

As it was mentioned before, likelihood is being built using MVA. The cut is chosen at the value when the significance  $\frac{S}{\sqrt{S+B}}$  reaches a maximum.

# Description of the flavor-tag consistency method

Estimation of  
btagging  
efficiency

Igor Marfin

Outline

Calibration  
methods

Likelihood ratio  
technique

Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary

Within the SM, top quarks are expected to decay almost to W boson accompanied by a b-quark.

In the semimounic ttbar events, given b efficiency and non-b mistag rate, the number of events with  $n_b$  tagged b-jets and  $n_{nonb}$  tagged nonb-jets can be predicted from MC.

By enforcing a consistency between the predicted number of events with no,one,two and more tagged jets to the actual number of observed events with that particular combination, the b-tag and non-btag efficiencies can be measured.

- The preselection of events is done.
- The MVA selection is performed to suppress the remain background (see the previous method).
- The following log-likelihood  $L = -2 \log \prod_n P(N_n, \bar{N}_n)$  must be minimized. Here  $N_n, \bar{N}_n, P$  are the measured number of events with  $n = 0, 1, 2$  tagged jets, the expected number of events, the Poisson distribution.
- The function  $\chi^2 = \sum_n \frac{(N_n - \bar{N}_n)^2}{\bar{N}_n}$  is minimized instead of the log-likelihood function.





- The expected number of events  $\bar{N}_n$  is calculated as
 
$$\bar{N}_n = L \times \sigma_{t\bar{t}b\bar{a}r} \times \epsilon_{sel}^{t\bar{t}b\bar{a}r} \times \sum_{i,j} f_{ij}^{t\bar{t}b\bar{a}r} \times \sum_{i+j=n}^{i \leq i, j \leq j} [C_i^i \epsilon_b^i \times (1 - \epsilon_b)^{(i-i)} \times C_j^j \epsilon_{nonb}^j \times (1 - \epsilon_{nonb})^{(j-j)}] +$$

$$+ L \times \sigma_{bkg} \times \epsilon_{sel}^{bkg} \times \sum_{i,j} f_{ij}^{bkg} \times \sum_{i+j=n}^{i \leq i, j \leq j} [\dots],$$
- where  $L, \sigma_{t\bar{t}b\bar{a}r(bkg)}, \epsilon_{sel}^{t\bar{t}b\bar{a}r(bkg)}$  are the luminosity, cross section of signal (background), the preselection and MVA combined efficiency.
- The coefficients  $f_{ij}^{t\bar{t}b\bar{a}r(bkg)}, C_i^i$  are the fraction of events with  $i, j$  of  $b-$  and  $nonb-$  jets respectively, and the binomial coefficients.
- The method gives  $\epsilon_b$  and  $\epsilon_{nonb}$

This method is based on the *TKinFitter* package . More details are available in **CMS AN 2005/025**.

- 4 constraints are used:  $mW_{lep}$ ,  $mW_{had}$ ,  $mTop_{lep}$ ,  $mTop_{had}$ .
- The constraints are Gaussian smeared:  $mW = 80.4\text{GeV} \pm \Gamma_W(2.1\text{GeV})$ ,  $mTop = 173.\text{GeV} \pm \Gamma_{Top}(12.7\text{GeV})$
- The parametrization is :  $\vec{p} = (E_T \cos\phi, E_T \sin\phi, E_T \sinh\eta)$ ,  $E = E_T \cosh\eta$
- Up to 7 jets are used to construct the  $\chi^2$ .
- Only such combination of 4 fitted jets with fit. muon and MET is used which has minimal  $\chi^2$  and has been converged.

Fig: Mass distribution of the fitted W boson

The invariants mass of fitted Whad cand		$\chi^2 / \text{ndf}$	6.853 / 9
Constant	13.5 $\pm$ 2.1		
Mean	80.58 $\pm$ 0.32		
Sigma	2.376 $\pm$ 0.297		

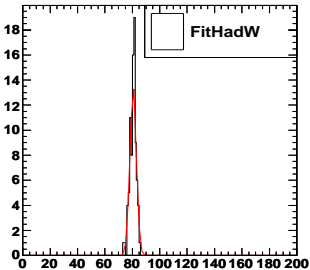
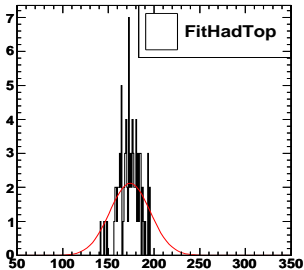


Fig: Mass distribution of the fitted top quark

Fitted tophad invariants mass		$\chi^2 / \text{ndf}$	17.71 / 34
Constant	2.121 $\pm$ 0.342		
Mean	173.8 $\pm$ 5.6		
Sigma	22.16 $\pm$ 8.95		





The selection of semimuonic ttbar events among overwhelming background is done using the following steps.

- SisCone algorithm with  $\Delta R = 0.5$  is used to construct jets.
- JES Corrections L2L3 are used
- The lepton impact parameter  $d_0$  is calculated with respect to the offline Beamspot
- $Reliso = (E_{calo}(Iso) + P_T(tracker, Iso))/P_T(\mu)$

**Table:** The Selection derived from TOP-09-003

Step	Description
<b>Step1</b>	$\geq 4$ jets with $Pt > 30\text{GeV}$ (corrected), $\eta < 2.4$
<b>Step2</b>	One GM muon with : $Pt > 30\text{GeV}$ , $\eta < 2.1$ , $N(hits) \geq 11$ , $d_0 < 200\mu$ , $\chi^2/ndf < 10$ , $Reliso < 0.05$
<b>Step3</b>	veto on electrons (no electrons which are GsfElectron, $\eta < 2.5$ , $Pt > 30\text{GeV}$ , $Reliso < 0.05$ )
<b>Step4</b>	Reconstructed MET. No cuts on MET are applied

Three different sets of the kinematical variables were used to train MVA. Two of them consider properties of fitted objects.

- Kinematical variables from CMS NOTE 2006/013. KinFitter is used.
- Kinematical variables of TopEventSelection package with Kinematical Fit.
- Kinematical variables of TopEventSelection package without Kinematical Fit.

The first set of the variables was chosen at the same way as in CMS NOTE 2006/013

## Kinematical variables from CMS NOTE 2006/013

- $pT_{hadtop}, \eta_{hadtop}$
- $pT_{lepton}, \eta_{lepton}$
- $pT_{hadB}, \eta_{hadB}$
- $pT_{lepB}, \eta_{lepB}$
- $\Delta\phi(hadB, hadtop), \Delta\theta(hadQ, hadQBar)$
- $\Delta\phi(hadB, hadW), \Delta\phi(lepB, lepW)$
- $pT_{3jet}/pT_{4jet}$
- $\Delta M(lepton, lepW)$
- $\Delta R(lepton, lepW)$
- $\Delta M(hadtop, hadW)$
- $\Delta R(hadtop, hadW)$
- $Prob(\chi^2)$

The second set of the kinematical variables was build on fitted objects as well as non-fitted ones

## Kinematical variables of TopEventSelection package with/without Kinematical Fit

- $sum_{E_T} = \sum_{i=1}^4 E_T(j_i)$
- $relEt1 = E_T(j_1)/sum_{E_T}$
- $MET.Et()$
- $mindijetmass = Min(Mass(j_i, j_k))/\sum_k M(j_k)$
- $maxdijetmass = Max(Mass(j_i, j_k))/\sum_k M(j_k)$
- $mindRjetlepton = Min(\Delta R(muon, j_k))$
- $lepeta = abs(\eta(muon))$
- $dphiMETlepton = \Delta\phi(MET, muon)$

## Procedure of the kinematical fit

- Perform non-linear least-square kinematical fit.
- The fit produces combinations of four fitted jets corresponding to  $b$ -, light-quarks fitted muon and neutrino.
- Choose the combination with minimal  $\chi^2$ .
- Construct  $W_{lep}$ ,  $W_{had}$ ,  $t_{lep}$ ,  $t_{had}$  candidates from fitted objects.
- Use kinematical variables to train MVA.

- PYTHIA6 from SUMMER09@7TeV samples.
- <https://twiki.cern.ch/twiki/bin/view/CMS/ProductionSummer2009at7TeV>
- ttbar events: /TTbar/Summer09-MC\_31X\_V3\_7TeV-v5/GEN-SIM-RECO/
- W+jets events:  
/Wmunu/Summer09-MC\_31X\_V3\_7TeV-v1/GEN-SIM-RECO
- QCD events:  
/InclusiveMu15\_Pt30/Summer09-MC\_31X\_V3\_7TeV-v1/GEN-SIM-RECO
- Wbb :/Wbb0Jets-alpgen/Summer09-MC\_31X\_V3\_7TeV-v2/GEN-SIM-REC
- Zbb: /Zbb0Jets-alpgen/Summer09-MC\_31X\_V3\_7TeV-v1/GEN-SIM-RECO

All events from datasets were asked to process.

*'total\_number\_of\_events = -1'*

# Preselecton. Cut Flow Table

Kinematical variables of TopEventSelection package with KinFit  
Table for TTbarPresel

Xsection in pb 94.3

Evnt tot	Evnt jets rej	Evnt muon rej	Evnt ele rej	Evnt rej MET	GenEvt rej
209630	191503	9278	8963	8963	8963

Table for WjetsPresel

Xsection in pb 7899

Evnt tot	Evnt jets rej	Evnt muon rej	Evnt ele rej	Evnt rej MET	GenEvt rej
2022023	790086	283	283	283	283

Table for WbbPresel

Xsection in pb 5.0724

Evnt tot	Evnt jets rej	Evnt muon rej	Evnt ele rej	Evnt rej MET	GenEvt rej
377004	306843	22	21	21	21

Estimation of  
btagging  
efficiency

Igor Marfin

Outline

Calibration  
methods

Likelihood ratio  
technique

Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary



Kinematical variables of TopEventSelection package with KinFit  
Table for ZbbPresel  
Xsection in pb 1.8046

Evnt tot	Evnt jets rej	Evnt muon rej	Evnt ele rej	Evnt rej MET	GenEvt rej
36618	30506	6	6	6	6

Table for QCDDPresel  
Xsection in pb 6.116e+07

Evnt tot	Evnt jets rej	Evnt muon rej	Evnt ele rej	Evnt rej MET	GenEvt rej
6411818	6409883	3	3	3	3





# MVA Training. Cut Flow Table

Estimation of  
btagging  
efficiency

Igor Marfin

Outline

Calibration  
methods

Likelihood ratio  
technique  
Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary

Kinematical variables of TopEventSelection package with KinFit

Table for TTbarTrainMVA

Xsection in pb 94.3

Evnt tot	Accept to train
4696	4414

Table for WjetsTrainMVA

Xsection in pb 7899

Evnt tot	Accept to train
283	265

Table for WbbTrainMVA

Xsection in pb 5.1528

Evnt tot	Accept to train
21	21



Kinematical variables of TopEventSelection package with KinFit  
Table for ZbbTrainMVA  
Xsection in pb 1.8046

----- -----
Evt tot   Accept to train
----- -----
6   6
----- -----

Table for QCDTrainMVA  
Xsection in pb 6.116e+07

----- -----
Evt tot   Accept to train
----- -----
3   3
----- -----



# MVA Efficiencies

Fig:TopEventSelection kinematical variables with  
kinfit

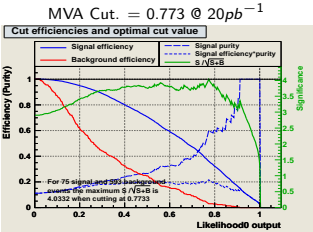


Fig:TopEventSelection kinematical variables  
w/o kinfit

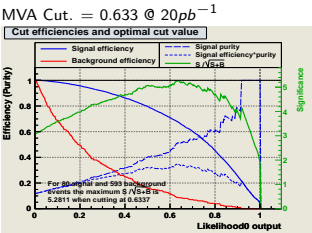


Fig:CMS NOTE 2006/013 kinematical variables  
with kinfit

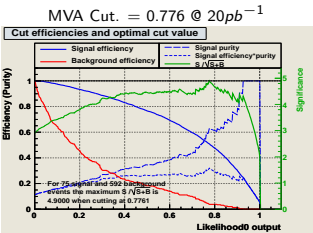


Table: Sig and Bkg events accepted to train  
 $@ 20pb^{-1}$

Kin. set	Sig	Bkg
TQAF +kinfit	75	593
TQAF -kinfit	80	594
CMS NOTE	75	592

# Normalized distributions of the variables. I

Estimation of  
btagging  
efficiency

Igor Marfin

Outline

Calibration  
methods

Likelihood ratio  
technique  
Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary

Fig: **TopEventSelection** +kinfit kinematical

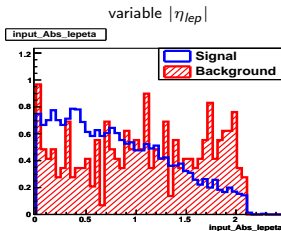


Fig: **CMS NOTE 2006/013** kinematical variable

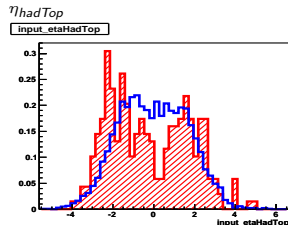


Fig: **TopEventSelection** -kinfit kinematical

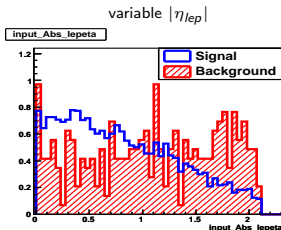
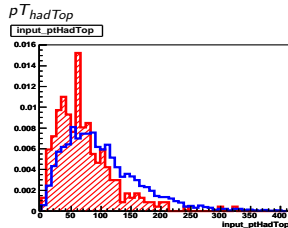


Fig: **CMS NOTE 2006/013** kinematical variable



# Normalized distributions of the variables. II

Estimation of  
tagging  
efficiency

Igor Marfin

Outline

Calibration  
methods

Likelihood ratio  
technique  
Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
tagging eff.

Summary

Fig: **TopEventSelection** +kinfit kinematical  
variable  $\max(M_{j1,j2})$

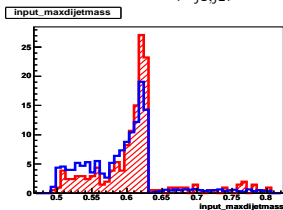


Fig: **TopEventSelection** -kinfit kinematical  
variable  $\max(M_{j1,j2})$

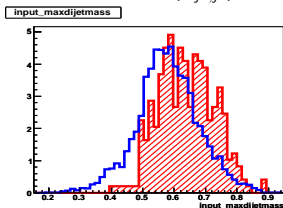


Fig: **CMS NOTE 2006/013** kinematical variable  
 $\Delta(\theta(q), \theta(\bar{q}))$

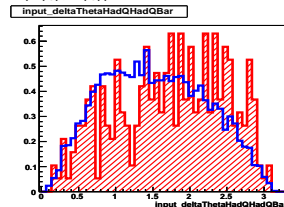
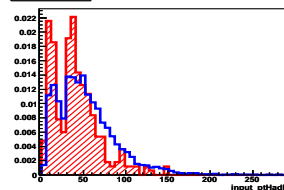


Fig: **CMS NOTE 2006/013** kinematical variable

$pT_{hadB}$   
**input\_ptHadB**



One needs to calculate  $x_b$  fraction and  $f_{ij}$  coefficients from MC events. They will be used for the estimation of btagging efficiency from real data later.

Kinematical variables of TopEventSelection package with KinFit  
Table for TTbarAnalEstimMVA  
Xsection in pb 9.430000e+01  
Lumi: 2.000000e+01

Evnt tot	Evnt pass pres	Presel Eff	Evnt pass mva	MVA Eff
1886	80.6384	0.04275	27.4061	0.33986

<x_b>	btag-j/b-j	nonbtag/nonb
0.3087	0.79551	0.10332



Kinematical variables of TopEventSelection package with KinFit  
Table for WjetsAnalEstimMVA  
Xsection in pb 7.899000e+03  
Lumi: 2.000000e+01

Evnt tot	Evnt pass pres	Presel Eff	Evnt pass mva	MVA Eff
157980	22.1107	0.00013	2.34389	0.106
<x_b>	btag-j/b-j	nonbtag/nonb		
0.0124	0	0.06918		

Kinematical variables of CMS NOTE  
Table for WjetsAnalEstimMVA  
Xsection in pb 7.899000e+03  
Lumi: 2.000000e+01

Evnt tot	Evnt pass pres	Presel Eff	Evnt pass mva	MVA Eff
157980	22.1107	0.00013	0.46877	0.0212
<x_b>	btag-j/b-j	nonbtag/nonb		
0	0	0.03333		



Kinematical variables of CMS NOTE  
Table for TTbarAnalEstimMVA  
Xsection in pb 9.430000e+01  
Lumi: 2.000000e+01

Evnt tot	Evnt pass pres	Presel Eff	Evnt pass mva	MVA Eff
1886	80.6384	0.04275	38.3788	0.4759
<x_b>	btag-j/b-j	nonbtag/nonb		
0.3034	0.79574	0.09876		

Kinematical variables of TopEventSelection package w/o KinFit  
Table for TTbarAnalEstimMVA  
Xsection in pb 9.430000e+01  
Lumi: 2.000000e+01

Evnt tot	Evnt pass pres	Presel Eff	Evnt pass mva	MVA Eff
1886	80.6384	0.04275	54.4858	0.67568
<x_b>	btag-j/b-j	nonbtag/nonb		
0.3067	0.79548	0.1014		





Table: Comparison of different kinematical sets.

Kin. set	Sig. presel.	$\frac{S}{S+B}$ presel.	Sig. MVA	$\frac{S}{S+B}$ MVA
TQAF + kinfit	80.6384	0.11	27.4061	0.9209
CMS NOTE + kinfit	80.6384	0.11	38.3788	0.9876
TQAF -kinfit	80.6384	0.11	54.4858	0.9679



# btagging efficiency and $x_b$ distribution plots

There are several plots corresponded to 'trackCountingHighEffBJetTags'  
btagging algo at the 'loose' operation point.

Fig:TopEventSelection +kinfit **btagging eff.**

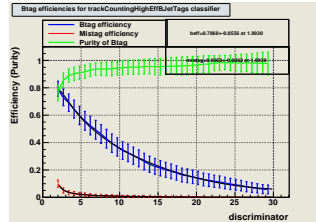


Fig:CMS NOTE 2006/013 **btagging eff.**

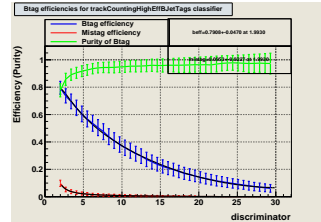


Fig:TopEventSelection +kinfit  $x_b$

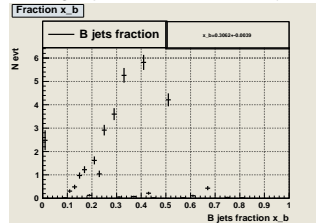
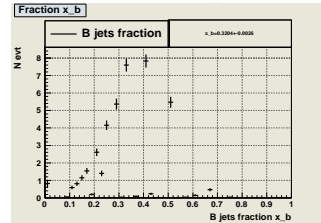


Fig:CMS NOTE 2006/013  $x_b$



# Estimation of btagging efficiency from pseudo-experiments

Outline

Calibration  
methods

Likelihood ratio  
technique  
Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary



- 'trackCountingHighEffBJetTags' btagging algo at the 'loose' operation
- 300 pseudo-experiments on  $20pb^{-1}$  data.
- TopEventSelection & kinfit kinematical variables are used.

Fig:Likelihood technique

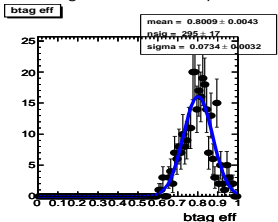


Fig:Flavor-tag consistency method

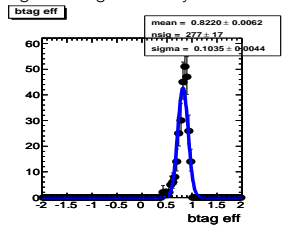


Fig:Likelihood technique. Pull distribution

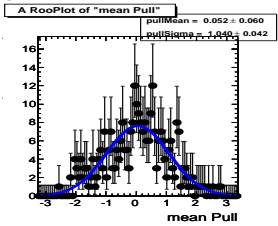
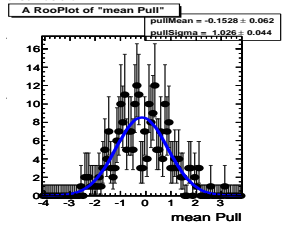


Fig:Flavor-tag consistency method. Pull distribution





■ 300 pseudo-experiments on  $100\text{pb}^{-1}$  data.

Fig:Likelihood technique

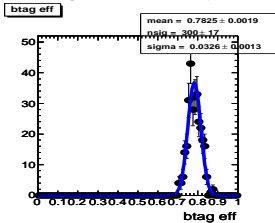


Fig:Likelihood technique. Pull distribution

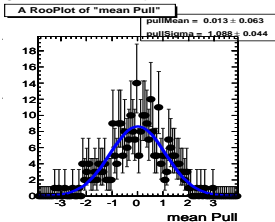


Fig:Flavor-tag consistency method

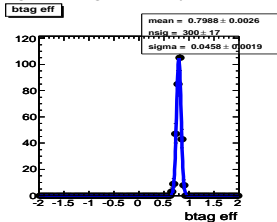
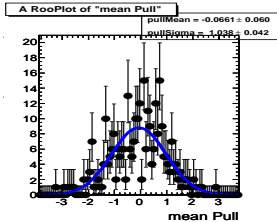


Fig:Flavor-tag consistency method. Pull distribution



**Table:** Estimation of the btagging efficiency. TopEventSelection & kinfit kinematical variables

Method	Lumi, $pb^{-1}$	estim. eff.	MC eff.
<b>Likelihood technique</b>	20	$0.80 \pm 0.07$	$0.78 \pm 0.05$
<b>Flavor-tag consistency method</b>	20	$0.82 \pm 0.10$	$0.78 \pm 0.05$
<b>Likelihood technique</b>	100	$0.78 \pm 0.03$	$0.78 \pm 0.02$
<b>Flavor-tag consistency method</b>	100	$0.79 \pm 0.04$	$0.78 \pm 0.02$

- 'trackCountingHighEffBJetTags' btagging algo at the 'loose' operation
- 300 pseudo-experiments on  $20\text{pb}^{-1}$  data.
- CMS NOTE kinematical variables are used.

Fig:Likelihood technique

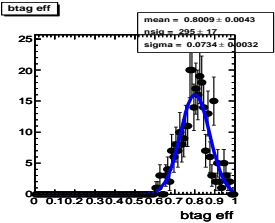


Fig:Flavor-tag consistency method

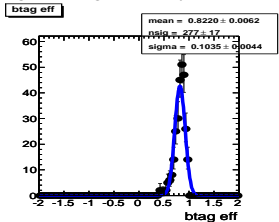


Fig:Likelihood technique. Pull distribution

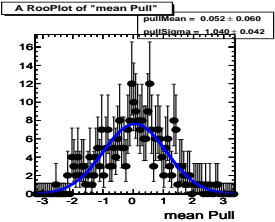
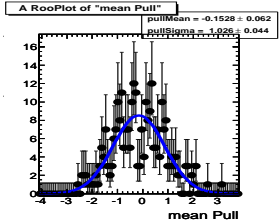


Fig:Flavor-tag consistency method. Pull distribution





■ 300 pseudo-experiments on  $100\text{pb}^{-1}$  data.

Fig:Likelihood technique

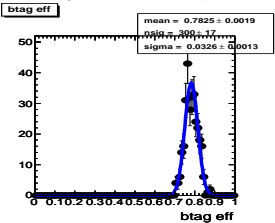


Fig:Flavor-tag consistency method

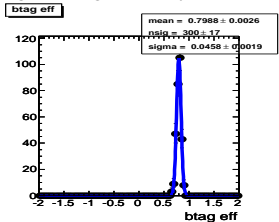


Fig:Likelihood technique. Pull distribution

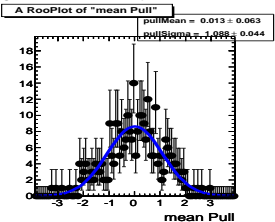
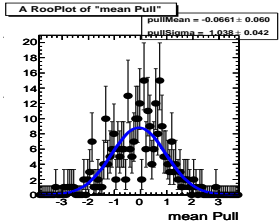


Fig:Flavor-tag consistency method. Pull distribution



**Table:** Estimation of the btagging efficiency. CMS NOTE kinematical variables

Method	Lumi, $pb^{-1}$	estim. eff.	MC eff.
<b>Likelihood technique</b>	20	$0.80 \pm 0.07$	$0.78 \pm 0.05$
<b>Flavor-tag consistency method</b>	20	$0.82 \pm 0.10$	$0.78 \pm 0.05$
<b>Likelihood technique</b>	100	$0.78 \pm 0.03$	$0.78 \pm 0.02$
<b>Flavor-tag consistency method</b>	100	$0.79 \pm 0.04$	$0.78 \pm 0.02$

**Table:** Estimation of the btagging efficiency. TopEventSelection w/o kinfit kinematical variables

Method	Lumi, $pb^{-1}$	estim. eff.	MC eff.
<b>Likelihood technique</b>	20	$0.77 \pm 0.06$	$0.78 \pm 0.05$
<b>Flavor-tag consistency method</b>	20	$0.83 \pm 0.08$	$0.78 \pm 0.05$
<b>Likelihood technique</b>	100	$0.78 \pm 0.02$	$0.78 \pm 0.02$
<b>Flavor-tag consistency method</b>	100	$0.80 \pm 0.03$	$0.78 \pm 0.02$



Outline

Calibration  
methods

Likelihood ratio  
technique

Flavor-tag  
consistency  
method

Kinematical  
Fit

Events  
preselection

Kinematical  
variables

Samples used  
for study

MVA Training

MVA Analysis

Estimation of  
btagging eff.

Summary

- The study has been finished. The following was done:
- MVA selection of the process have been developed and implemented in CMSSW.
- Two methods of btagging calibration were studied.
- Three sets of kinematic variables were used.
- The kinematic variables determined as in CMS NOTE 2006/013 are the best choice to make the estimation of btagging efficiency.
- Thanks for your attention.

