

1 Multijet background modelling with the Matrix Method

1.1 Motivation for a data-driven technique

The multijet background, also referred to as QCD background, is one of the most important at hadron colliders. It is also one of the most difficult to model since it is intimately interwoven with various detector-related quantities such as lepton identification and isolation criteria.

Since the semi-leptonic $t\bar{t}$ signal is comprised of 4 jets and a W boson, decaying into an lepton and a neutrino, an isolated lepton is required by the analysis cuts. In *muon + jets* topology, the QCD background is comprised of muons coming from semi-leptonic decays. Such a muon, if its parent jet is not well reconstructed or if it is emitted out of it, can be identified as isolated and thus wrongly accepted as the decay product of the W boson. This component is also present in the *electron + jets* topology. In the latter however, another effect must be taken into account, as pion-rich, electromagnetic-like jets can be wrongly identified as isolated electrons. The presence of this component makes the QCD background higher and more difficult to model in the electron channel, as the two contributions cannot be easily disentangled.

Since the QCD background has an instrumental component, it is difficult to model adequately with event generators. It is therefore relevant to use data-driven techniques such as the Matrix Method, as it has been successfully done at the Tevatron.

1.2 The Matrix Method: principle

The QCD background is estimated by modelling the kinematics of multijet events using real data as input. The normalisation and shape of the background can be extracted from data through the reweighting of individual data events. The Matrix Method technique provides means of derivating those weights.

The Matrix Method boils down to solving a set of two equations with two unknowns. Those equations are built from a two-stage event selection process. First a data sample is selected requiring the lepton to satisfy loose isolation criteria; this is thereafter referred to as the loose sample. Then a sub-sample is selected, requiring the lepton to satisfy tight isolation criteria. Those tight criteria are chosen so that they match the final, analysis-optimised isolation criteria. This sample is thereafter referred to as the tight sample. Given those two samples, the following system of equations can be constructed:

$$\begin{aligned} N_L &= N_l + N_{QCD} \\ N_T &= \epsilon N_l + f N_{QCD} \end{aligned} \tag{1}$$

where N_L (N_T) is the number of data events in the loose (tight) sample passing the selections, N_l is the number of signal-like leptons originating from the W decay and N_{QCD} the number of misidentified QCD leptons. The quantity ϵ , also referred to as signal efficiency, denotes the fraction of loose signal-like leptons subsequently passing the tight cut; f , also referred to as fake rate, denotes the fraction of loose QCD leptons passing the tight cut.

Since the tight selection criteria are the same as those used in the final analysis, this system of equations is solved for $N_{QCD}^T = f N_{QCD}$:

$$N_{QCD}^T = \frac{f}{\epsilon - f} (\epsilon N_L - N_T) \tag{2}$$

This result however cannot be used directly as the parameters ϵ and f can depend on event kinematics \vec{k} and have to be modelled accordingly. In order to take those dependencies into account, an unbinned method is used in which a weight is constructed. With this method, all the events are used in the background estimation, as the event weight is calculated differently for a loose event and for a tight event. The expression of the weight w_i is

$$w_i = \frac{f(\vec{k}_i)}{\epsilon(\vec{k}_i) - f(\vec{k}_i)} (\epsilon(\vec{k}_i) - T_i) \quad (3)$$

where $T_i = 1$ if event i satisfies both the loose and tight criteria and $T_i = 0$ if event i satisfies only the loose criteria. The QCD prediction N_{QCD}^T in the final sample is therefore the sum of weights calculated over the loose sample:

$$N_{QCD}^T = \sum_{i=1}^{N_L} w_i \quad (4)$$

The two parameters f and ϵ must be measured on independent and relevant samples, which can in principle be selected from collision data.

Both loose and tight selections are defined by applying the analysis selection cuts and varying the lepton isolation criteria. In the electron channel, the loose sample is defined by requiring $\Delta R(\text{electron, closest jet}) > 0.3$; the tight sample is defined by applying a relative isolation cut $reliso(\text{electron}) < 0.1$ in addition to the ΔR cut. In the muon channel the loose sample is defined by requiring $\Delta R(\text{muon, closest jet}) > 0.3$ and $reliso(\text{muon}) < 0.1$; the tight sample is defined by tightening the relative isolation cut to $reliso(\text{muon}) < 0.05$. In both cases the tight sample is therefore a tighter sub-sample of the loose sample.

1.3 Estimation of fake rate and signal efficiency

1.3.1 Fake rate

The fake rate f can be interpreted as the probability for a mis-identified, QCD-like lepton passing the loose selection cut to subsequently pass the tight selection cut. It can be measured in collision data by applying the two-stage loose / tight selection on a sample so selected as to be as QCD-like as possible. This QCD-like sample must be as independent as possible from the final analysis sample. An anti-isolation cut, which inverts the lepton isolation criteria and selects events with non-isolated lepton, results in having the fake rate directly depend on the isolation criteria chosen in the analysis. This technique is therefore not used.

Since the semi-leptonic $t\bar{t}$ signal is a $W + jets$ -like signature containing real \cancel{E}_T coming from one W boson decaying into a lepton and a neutrino, missing transvers energy can be used as a handle on the QCD background. Indeed, signal events contain a neutrino coming from a W decay and therefore feature relatively high \cancel{E}_T while QCD-like events contain no neutrino, or neutrinos coming from semi-leptonic decays in jets, and therefore features relatively low \cancel{E}_T . This difference is exploited to model the fake rate. The analysis cuts are applied on data, then a \cancel{E}_T cut is added. In the muon channel, the chosen cut is $\cancel{E}_T < 5\text{GeV}$; in the electron channel the cut is chosen to be $\cancel{E}_T < 12\text{GeV}$. The \cancel{E}_T cut is not applied in the final analysis. Because we are interested in the QCD prediction in the semi-leptonic $t\bar{t}$ topology, it is important that all the analysis cuts, including that on the number of jets required in the event, be applied to the QCD-like sample. In the muon channel however, the fake rate is found to be statistically compatible when estimated in an $n_{jets} \geq 2$ sample with that estimated in an $n_{jets} \geq 4$ sample. In order to maximise the available statistics we therefore use the value estimated in the $n_{jets} \geq 2$ bin. In the electron channel mis-identified jets also contribute to the QCD background, making the fake rate more strongly dependent on the number of jets present in the event. We therefore use the value estimated in the $n_{jets} \geq 4$ bin, which is the bin we study in the analysis. The triangle cut (see section 1.6) is not applied when measuring the fake rate since it is designed to suppress QCD-like events.

Since the main isolation criterion we use is a cut on $reliso$, which by construction is a function of lepton p_T , we derive both fake rate and efficiency as a function of lepton p_T . Those quantities are therefore taken as the tight-to-loose ratio of the lepton p_T ; a fit is then performed and the best fit function use as a parametrisation.

Because contamination of real $W/Z + jets$ events can occur, even at low \cancel{E}_T , we remove this contribution using simulated $W + jets$ and Drell-Yann samples. We apply the same cuts on the simulated samples then subtract their contribution from the selected data sample, since the contamination can affect the final shape and normalisation of the QCD background

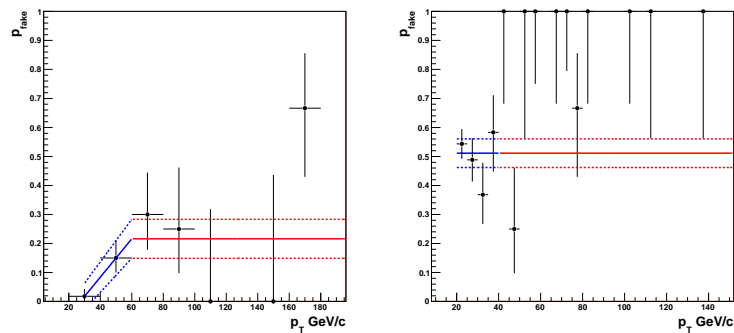
2010 Data Figure 1 shows the fit results for the electron and the muon channel. In the electron channel the fit function is chosen to be $p_0 + p_1 \cdot p_T^e$:

$$f^e(p_T^e) = -2.27889 \cdot 10^{-1} + 7.27134 \cdot 10^{-3} \cdot p_T^e \quad (5)$$

In the muon channel the fit function is chosen to be a constant:

$$f^\mu = 5.15437 \cdot 10^{-1} \quad (6)$$

Since statistics are low at high transverse momentum the fake rate is capped with the last value computed with acceptable statistics. This is indicated on the plots by a red line.



(a) Fake rate in the electron channel (b) Fake rate in the muon channel

Figure 1: Fake rate in (a) the electron channel and (b) the muon channel as a function of the lepton transverse momentum in 2010 data. The points are the loose-to-tight ratio of the lepton p_T in a QCD-like data sample; the overlaid blue curve is the best fit function found to describe those points, while the red line indicates the capping value. The functional forms are given by equations 5 and 6.

2011 data The same technique is applied in 2011 data; the resulting fits can be seen on figure 2. In the electron channel both fake rate and signal efficiency are derived as a function of p_T^e in two η_e bins corresponding to the barrel and endcap calorimeters; the fit gives

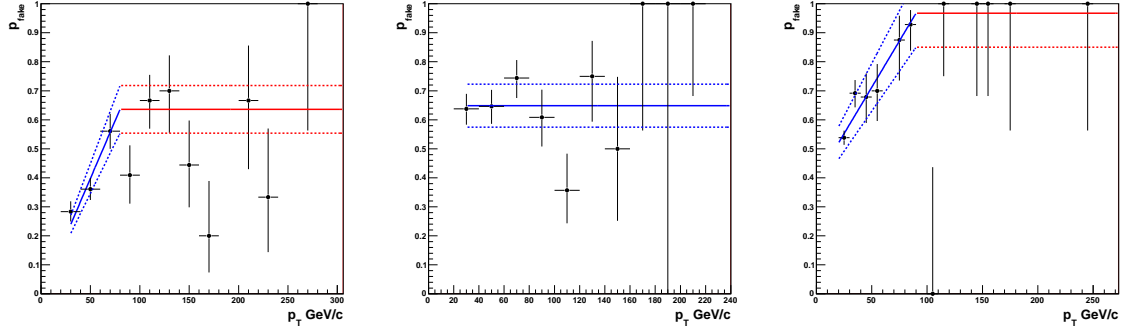
$$\begin{aligned} f^e(p_T^e, \eta_e < 1.45) &= 7.95157 \cdot 10^{-3} \cdot p_T^e \\ f^e(p_T^e, \eta_e \geq 1.45) &= 6.45072 \cdot 10^{-1} \end{aligned} \quad (7)$$

A linear function is chosen for the muon channel:

$$f^\mu = 3.92423 \cdot 10^{-1} + 6.38960 \cdot 10^{-3} * p_T^\mu \quad (8)$$

1.3.2 Signal efficiency

The signal efficiency ϵ can be interpreted as the probability for a signal-like lepton passing the loose selection cut to subsequently pass the tight selection cut. It can be measured in collision data by applying the two-stage loose / tight selection on a sample so selected as to be as signal-like as possible.



(a) Fake rate in the electron channel, $\eta_e < 1.45$ (b) Fake rate in the electron channel, $\eta_e \geq 1.45$ (c) Fake rate in the muon channel

Figure 2: Fake rate in (a) and (b) the electron channel and (c) the muon channel as a function of the lepton transverse momentum in 2011 data. The electron fake rate is measured in two bins of η_e . The points are the loose-to-tight ratio of the lepton p_T in a QCD-like data sample; the overlaid blue curve is the best fit function found to describe those points, while the red line indicates the capping value. The functional forms are given by equations 7 and 8.

This means that only events containing physically isolated leptons should be included in this sample. Ideally, a $Z \rightarrow ll$ sample, with $l = e, \mu$, should be selected from data and used to measure the signal efficiency. However, since the available integrated luminosity used in this analysis is relatively low, it is not possible to select such a sample with enough events to perform a reliable fit. $W \rightarrow e, \mu + jets$ simulation is therefore used for the electron and muon channel respectively. The analysis selections are applied on this sample, without any \cancel{E}_T cut.

2010 data Figure 3 shows the fit results for the electron and the muon channel respectively. In the electron channel the functional form $p_0 + p_1 \cdot \log(p_T^e)$ is used. The best parameters returned by the fit are

$$\epsilon^e(p_T^e) = 3.90864 \cdot 10^{-1} + 1.39888 \cdot 10^{-1} \cdot \log(p_T^e) \quad (9)$$

In the muon channel the fit function is of the form $p_0 + p_1 \cdot \log(p_T^\mu) + p_2 \cdot (p_T^\mu)^3$; the fit result is

$$\epsilon^\mu(p_T^\mu) = 5.52913 \cdot 10^{-1} + 9.88649 \cdot 10^{-2} \log(p_T^\mu) - 2.59743 \cdot 10^{-8} \cdot (p_T^\mu)^3 \quad (10)$$

2011 data We use the same technique in 2011 data. In the electron channel, the following functional form is used to model the signal efficiency:

$$\begin{aligned} \epsilon^e(p_T^e, \eta_e < 1.45) &= 3.46971 \cdot 10^{-1} + 1.36916 \cdot 10^{-1} \cdot \log(p_T^e) \\ \epsilon^e(p_T^e, \eta_e \geq 1.45) &= 7.45471 \cdot 10^{-1} + 5.14007 \cdot 10^{-2} \cdot \log(p_T^e) \end{aligned} \quad (11)$$

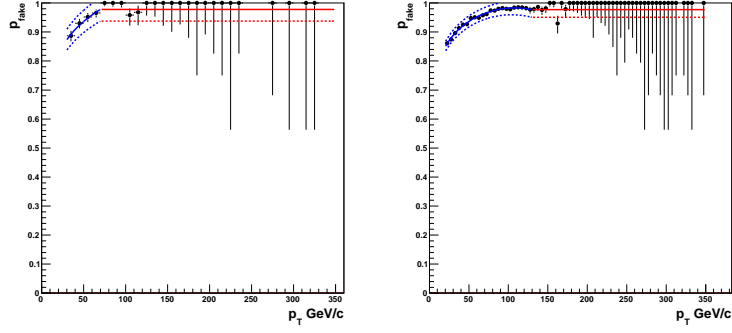
In the muon channel the following function is used:

$$\epsilon^\mu(p_T^\mu) = 2.11304 \cdot 10^{-1} + 1.71760 \cdot 10^{-1} \cdot \log(p_T^\mu) \quad (12)$$

Fake rate and signal efficiency for 2011 data can be seen on figure 4

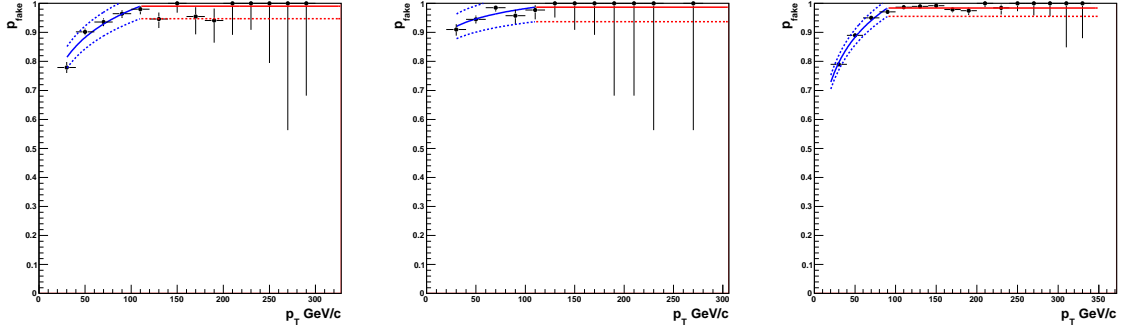
1.4 Normalisation and shape issues

muon channel In the muon channel the QCD background is mostly comprised of semi-leptonic decays in jets, and can therefore, in the first approximation, be said to feature only one component.



(a) Signal efficiency in the electron channel (b) Signal efficiency in the muon channel

Figure 3: Signal efficiency in (a) the electron channel and (b) the muon channel as a function of the lepton transverse momentum. The points are the loose-to-tight ratio of the lepton p_T in a $W \rightarrow e, \mu + jets$ simulated sample; the overlaid blue curve is the best fit function found to describe those points, while the red line indicates the capping value. The functional forms are given by equations 9 and 10.



(a) Signal efficiency in the electron channel, $\eta_e < 1.45$ (b) Signal efficiency in the electron channel, $\eta_e \geq 1.45$ (c) Signal efficiency in the muon channel

Figure 4: Signal efficiency in (a) and (b) the electron channel and (c) the muon channel as a function of the lepton transverse momentum. The electron fake rate is measured in two bins of η_e . The points are the loose-to-tight ratio of the lepton p_T in a $W \rightarrow e, \mu + jets$ simulated sample; the overlaid blue curve is the best fit function found to describe those points, while the red line indicates the capping value. The functional forms are given by equations 11 and 12.

The Matrix Method has already widely been shown, for instance at the Tevatron, to be excellent to model this type of background. We therefore use its output to predict both the shape and normalisation of the QCD background, the latter as an initial constraint for the final template fit. Figures 10 to ?? show a few important kinematic distributions before template fitting is performed; it can be seen that the prediction-to-data agreement, both in normalisation and shape, is good. The output of the Matrix Method is therefore considered to adequately describe the QCD background in the muon channel

electron channel In the electron channel the situation is more complicated because the QCD background has several components. In addition to semi-leptonic decays, which produce real but non-isolated electrons, other objects can fake an electron. This includes pion-rich jets but also photons. This has two direct consequences on the modelling of the QCD background:

1. Since there are two sources for background in the electron channel, the QCD background will be significantly higher than in the muon channel.
2. The gluon-gluon and photon+jets fractions of the background should ideally be disentangled, since those two topologies have different features: photon+jets events typically have less missing transverse energy, and the photon fake rate is *a priori* different from the jet fake rate.

Since the handle on QCD background in the Matrix Method is \cancel{E}_T , two problems can therefore arise from this last point if we use this technique the same way as in the muon channel:

1. The predicted shape can be incorrect; this however can be recovered by considering other dependencies for the fake rate and signal efficiency, i.e. deriving them as a function of the lepton p_T in bins of another adequately chosen variable; this is done by measuring the electron fake rate and signal efficiency in bins of the electron pseudo-rapidity. Doing so improves the shape of the electron η distribution; it has however very little impact on all the other kinematic distributions we investigated and does not significantly improve the shape agreement. In particular, the lepton sector is almost completely decoupled from the jet sector and therefore not improvement in the top mass distribution is observed. This has already been observed at the Tevatron.
2. Since the Matrix Method basically assumes a flat fake rate versus \cancel{E}_T , the fake rate must be extracted in a phase-space zone where both the gluon-gluon and the photon+jets regimes are similar; otherwise the Matrix Method will not normalise the predicted QCD background correctly.

We decide to use 2011 fake rate and signal efficiency on 2010 data so as to free ourselves from the very low statistics available in 2010 data; figure 5 shows that the 2010 and 2011 parameters are compatible within their error. Figures 14 to 17 show that the shape agreement thus obtained is very good.

1.5 Systematic uncertainties

The uncertainties on the fake rate and signal efficiency are given by varying the fit parameters by $\pm 1\sigma$. The varied fake rate and signal efficiency are then propagated to the analysis and used to calculate the Matrix Method weights. The resulting shapes are then used in the final template fit.

Figures 6 and 7 show the central QCD shape along with shapes obtained after shifting both fake rate and signal efficiency by $\pm 1\sigma$, for the hadronic top mass distribution, in the muon and the electron channel respectively. Those shapes are used in the final template fit from which we extract the production cross-section. Figures 8 and 9 show the deviation from the central shape when only one parameter, fake rate or signal efficiency, is shifted by $+1\sigma$, in the muon and electron channel respectively, in 2011 data.

–Yields to be added after output of template fit is known–

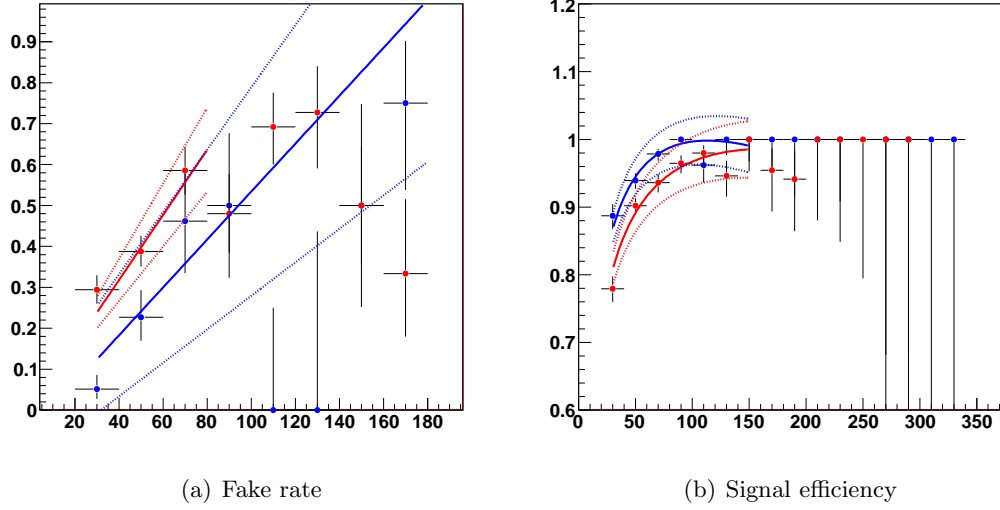


Figure 5: Comparison of (a) fake rate and (b) signal efficiency in the electron channel, as a function of the electron p_T , between the 2010 and 2011 epochs. This covers the whole pseudo-rapidity range as the 2010 statistics does not allow a break-down in several η regions. The points represent data while the curves show the output of the fit that is used in the analysis. The dotted line show the $\pm 1\sigma$ variation on those parameters. 2010 parameters are shown in blue while the 2011 ones are shown in red. 2010 and 2011 parameters are compatible within errors.

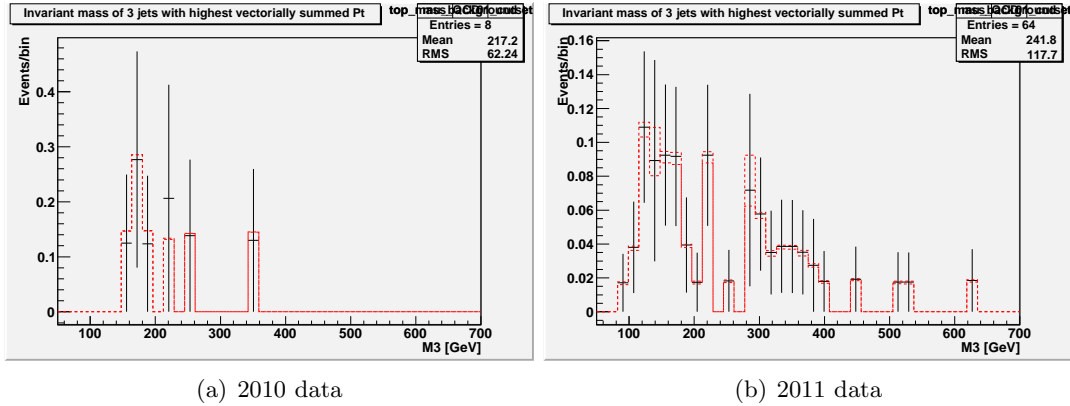
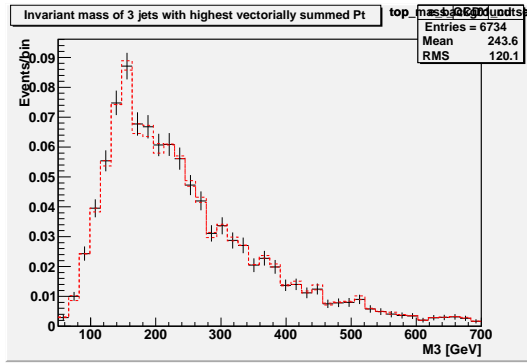


Figure 6: Shape of the QCD background as given by the Matrix Method for the hadronic top mass in the muon channel in (a) the 2010 data and (b) the 2011 data. The central value is shown in black while the shifted shapes are shown in the red, dotted-line histograms. Those shapes are used in the final template fit.



(a) 2011 data

Figure 7: Shape of the QCD background as given by the Matrix Method for the hadronic top mass in the electron channel in the 2011 data. The central value is shown in black while the shifted shapes are shown in the red, dotted-line histograms. Those shapes are used in the final template fit.

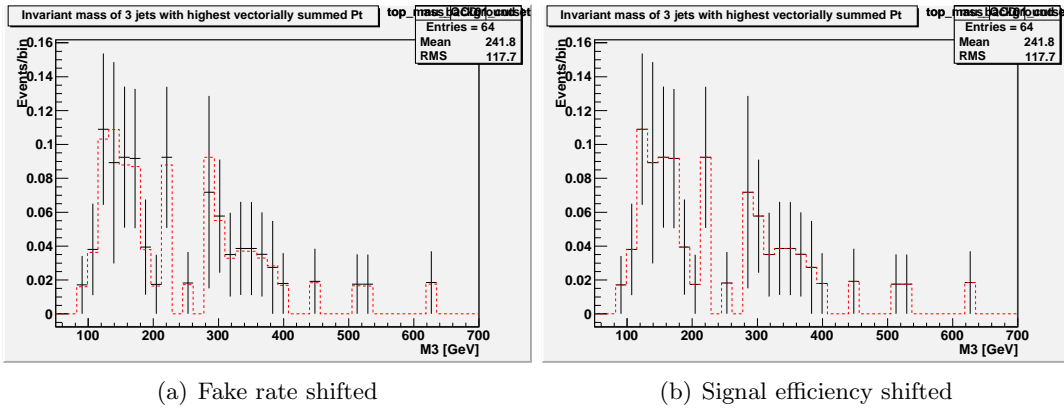


Figure 8: Shape of the QCD background as given by the Matrix Method for the hadronic top mass in the muon channel in the 2011 data. The overlaid, dotted curve shows the same shape when either (a) the fake rate or (b) the signal efficiency has been shifted by $+1\sigma$, the other parameter remaining constant.

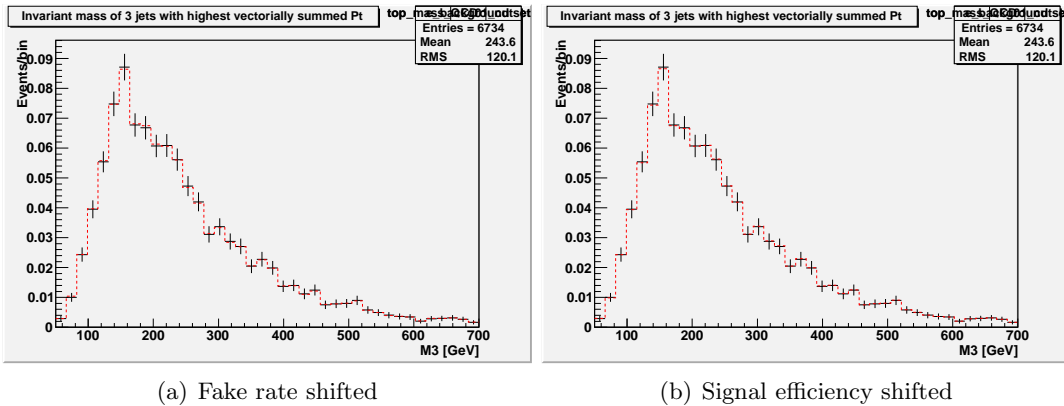


Figure 9: Shape of the QCD background as given by the Matrix Method for the hadronic top mass in the electron channel in the 2011 data. The overlaid, dotted curve shows the same shape when either (a) the fake rate or (b) the signal efficiency has been shifted by $+1\sigma$, the other parameter remaining constant.

1.6 Triangle cut

–This should not go in the QCD part of the note!– In order to improve agreement in the muon channel, we add an additional so-called triangle cut to the analysis selection. Since a good description of the data at low \cancel{E}_T , M_T^W and jet p_T is difficult to obtain in the muon channel in 2010 data, we exclude this region from the (M_T^W, \cancel{E}_T) 2-D plane by applying a triangle cut of $M_T^W > 115 - 0.5 \cancel{E}_T$. This cut is applied neither in the electron channel nor in 2011 data.

Figures 18 to 24 show distributions of the lepton transverse momentum, the missing transverse energy, the transverse mass of the leptonic W boson and the transverse momentum of the four leading jet in the muon channel before and after applying the triangle cut; they demonstrate the improvement in the overall agreement obtained by the application of the triangle cut.

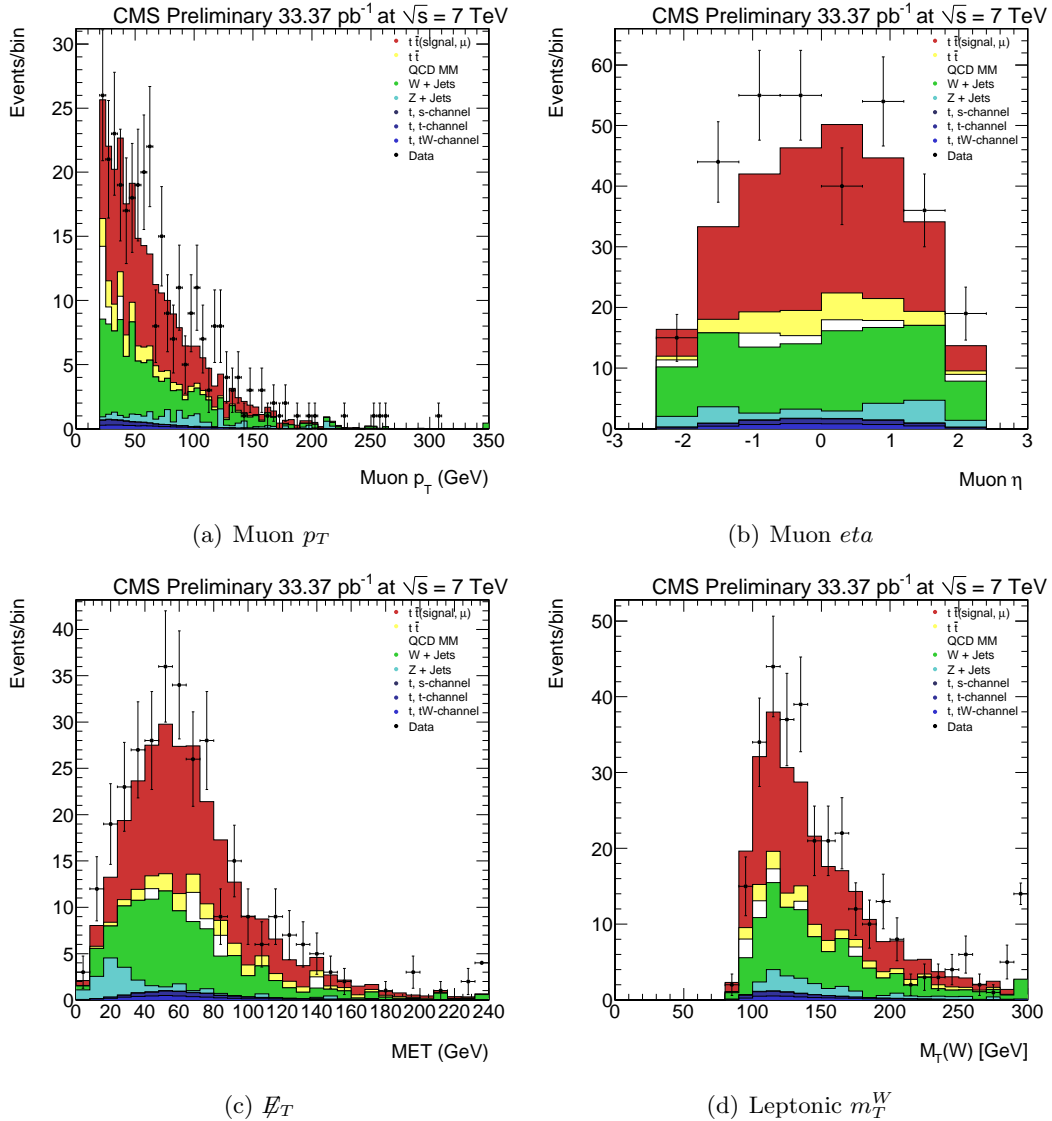


Figure 10: Leptonic kinematic distributions in the muon channel in 2010 data: (a) muon transverse momentum, (b) muon pseudo-rapidity, (c) missing transverse energy, (d) transverse mass of the leptonic W boson, after template renormalisation. There is good agreement between data and prediction.

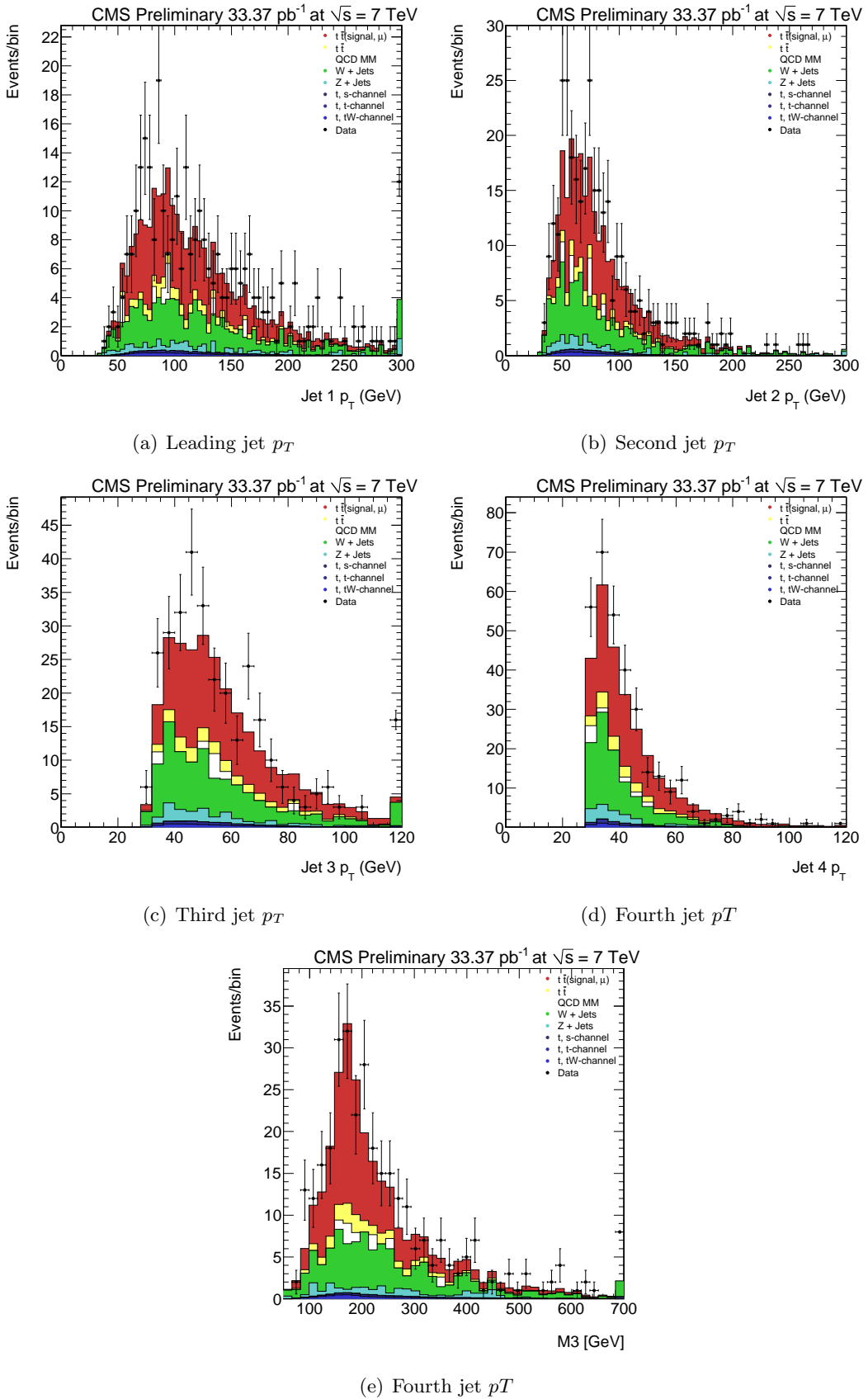


Figure 11: Jet kinematic distributions in the muon channel in 2010 data: (a) leading jet transverse momentum, (b) second-leading jet transverse momentum, (c) third-leading jet transverse momentum, (d) fourth-leading jet transverse momentum and (e) top mass, after template renormalisation. There is good agreement between data and prediction.

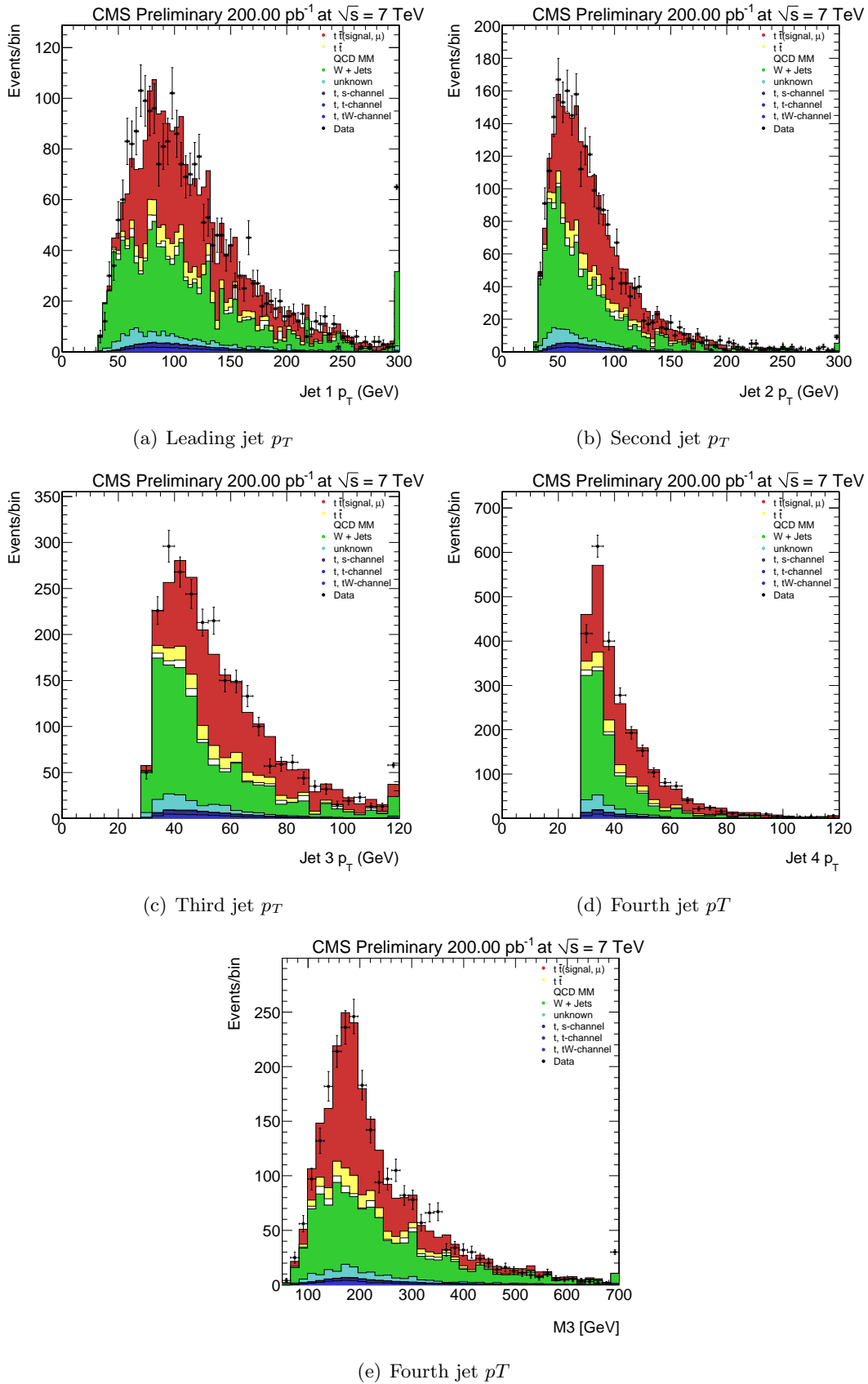


Figure 13: Jet kinematic distributions in the muon channel in 2011 data: (a) leading jet transverse momentum, (b) second-leading jet transverse momentum, (c) third-leading jet transverse momentum, (d) fourth-leading jet transverse momentum and (e) top mass, after template renormalisation. There is good agreement between data and prediction.

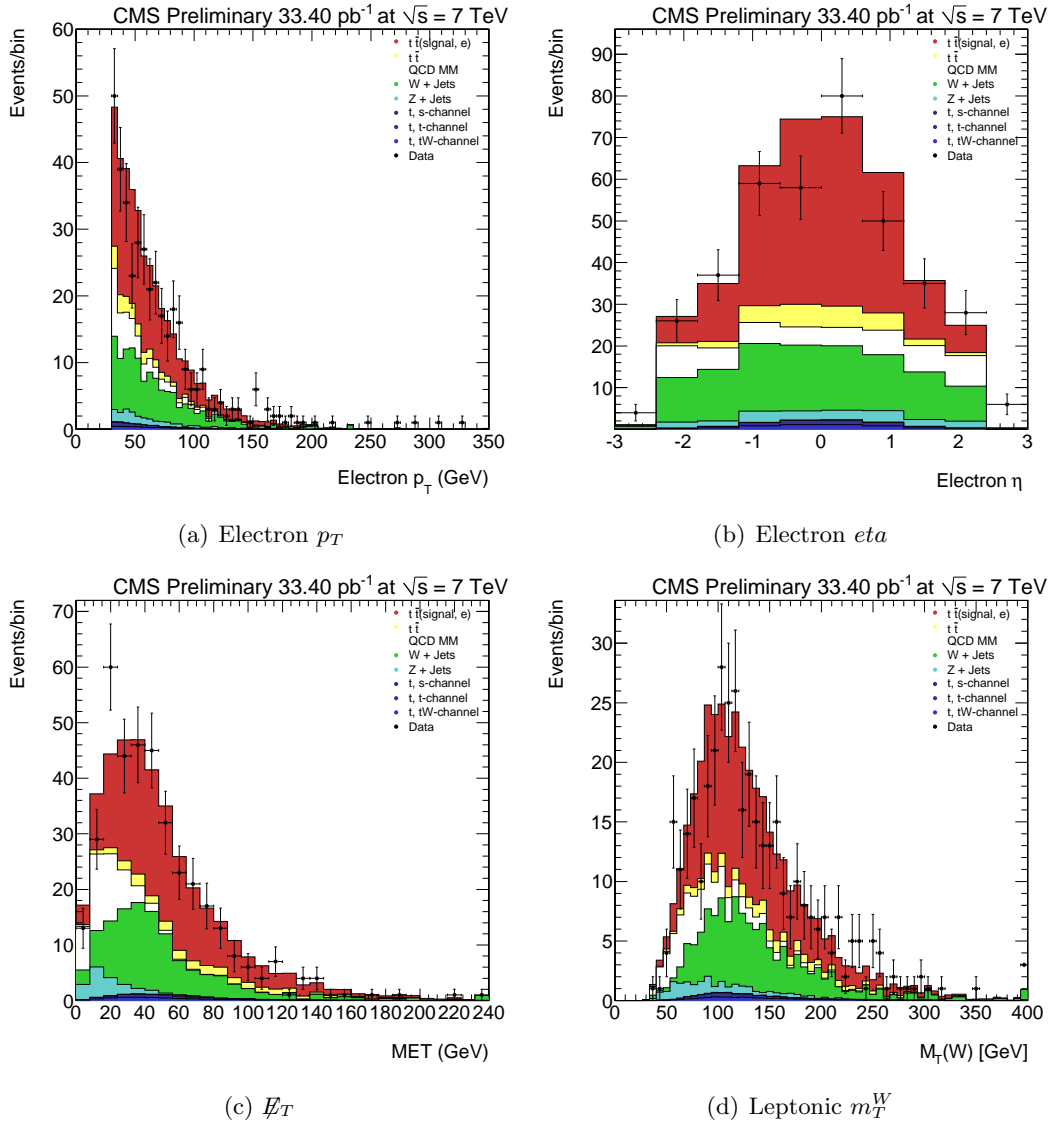


Figure 14: Leptonic kinematic distributions in the electron channel in 2010 data: (a) electron transverse momentum, (b) electron pseudo-rapidity, (c) missing transverse energy, (d) transverse mass of the leptonic W boson, after template renormalisation. There is good agreement between data and prediction.

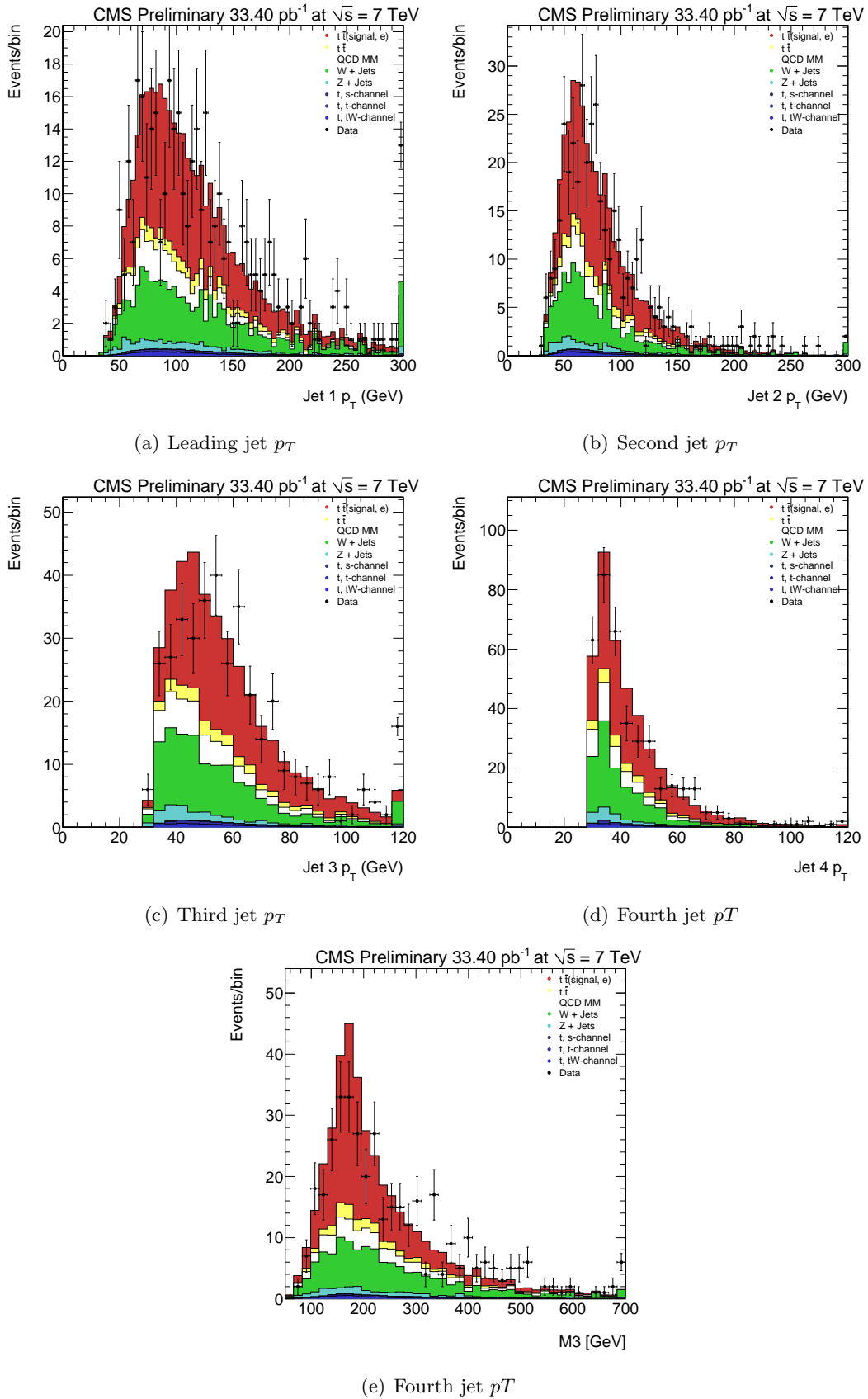


Figure 15: Jet kinematic distributions in the electron channel in 2010 data: (a) leading jet transverse momentum, (b) second-leading jet transverse momentum, (c) third-leading jet transverse momentum, (d) fourth-leading jet transverse momentum and (e) top mass, after template renormalisation. There is good agreement between data and prediction.

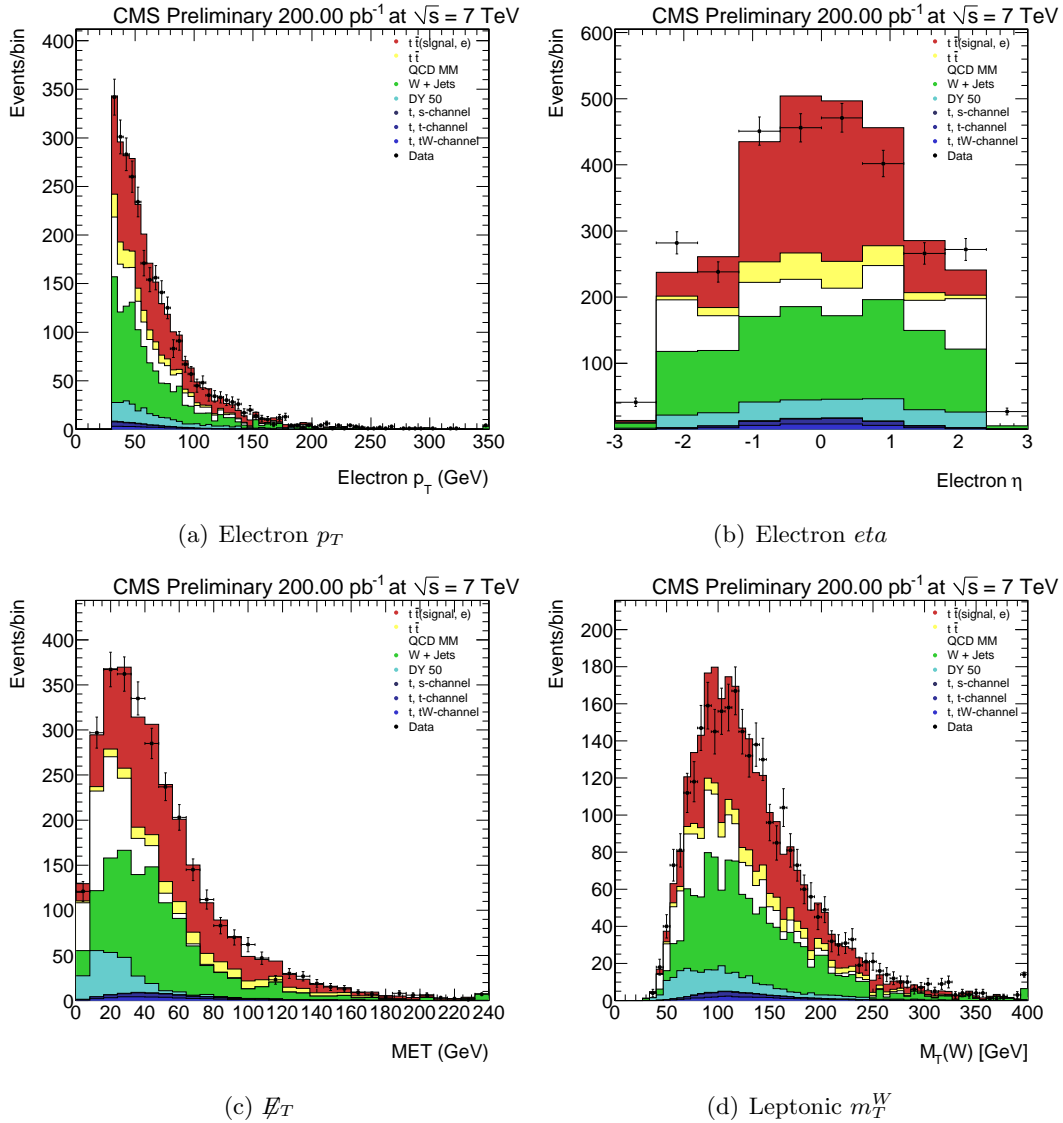


Figure 16: Leptonic kinematic distributions in the electron channel in 2011 data: (a) electron transverse momentum, (b) electron pseudo-rapidity, (c) missing transverse energy, (d) transverse mass of the leptonic W boson, after template renormalisation. There is good agreement between data and prediction.

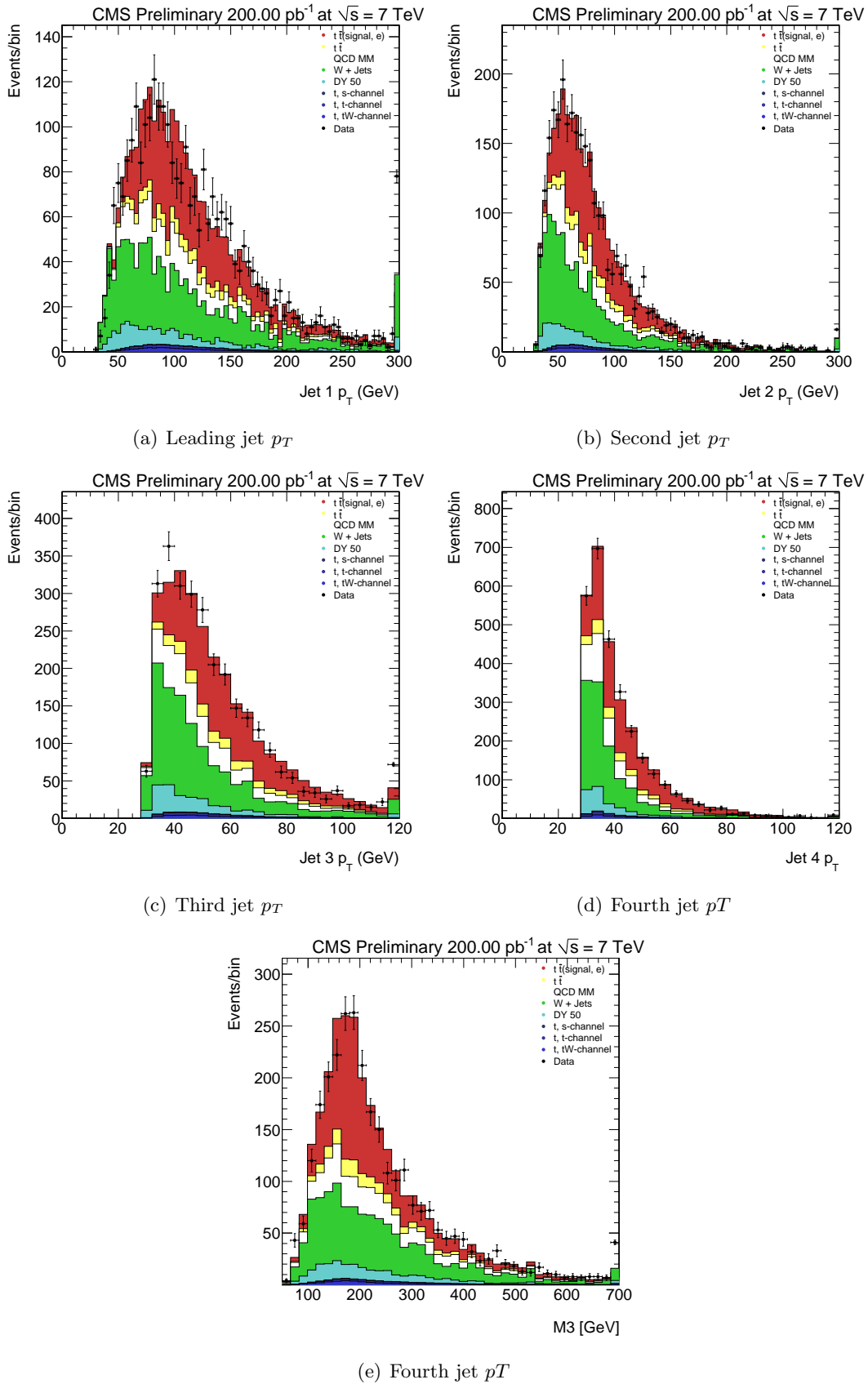


Figure 17: Jet kinematic distributions in the electron channel in 2011 data: (a) leading jet transverse momentum, (b) second-leading jet transverse momentum, (c) third-leading jet transverse momentum, (d) fourth-leading jet transverse momentum and (e) top mass, after template renormalisation. There is good agreement between data and prediction.

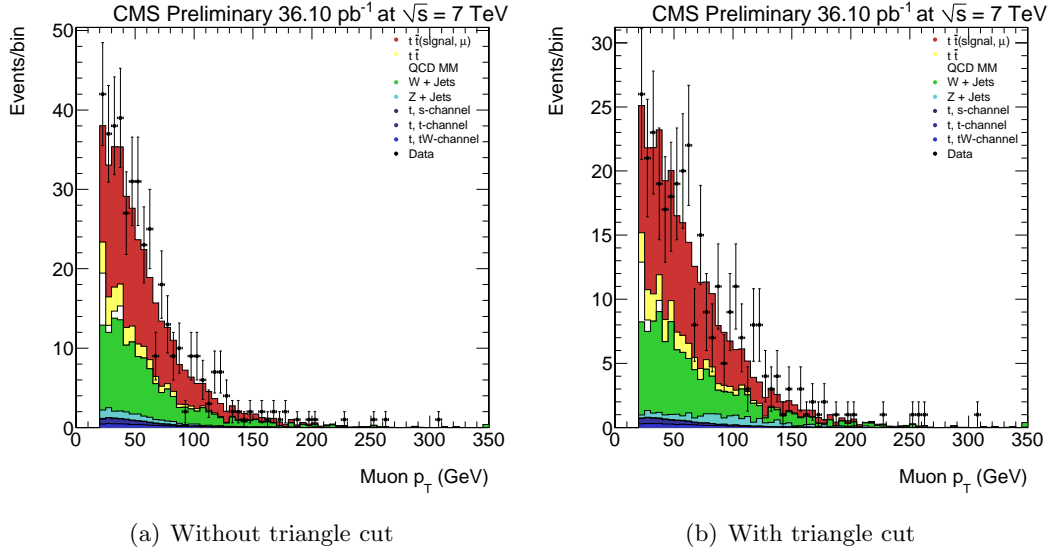


Figure 18: Distribution of the lepton transverse momentum in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

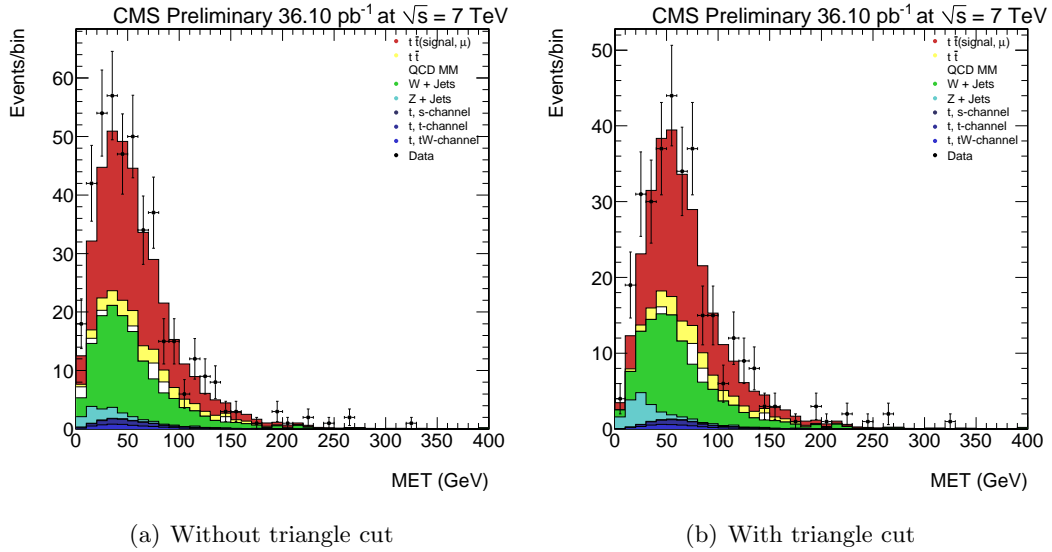


Figure 19: Distribution of the missing transverse energy in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

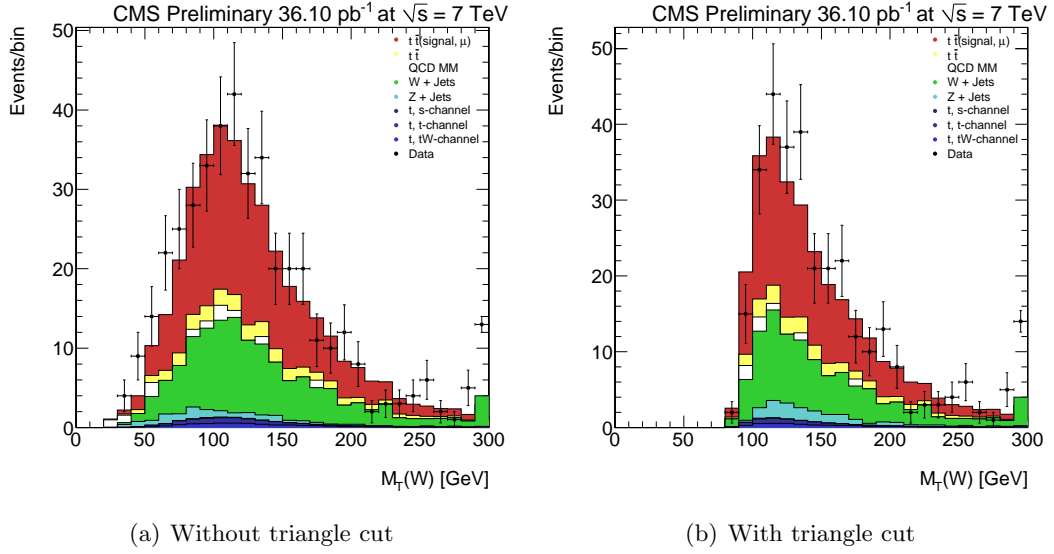


Figure 20: Distribution of the leptonic W boson transverse mass in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

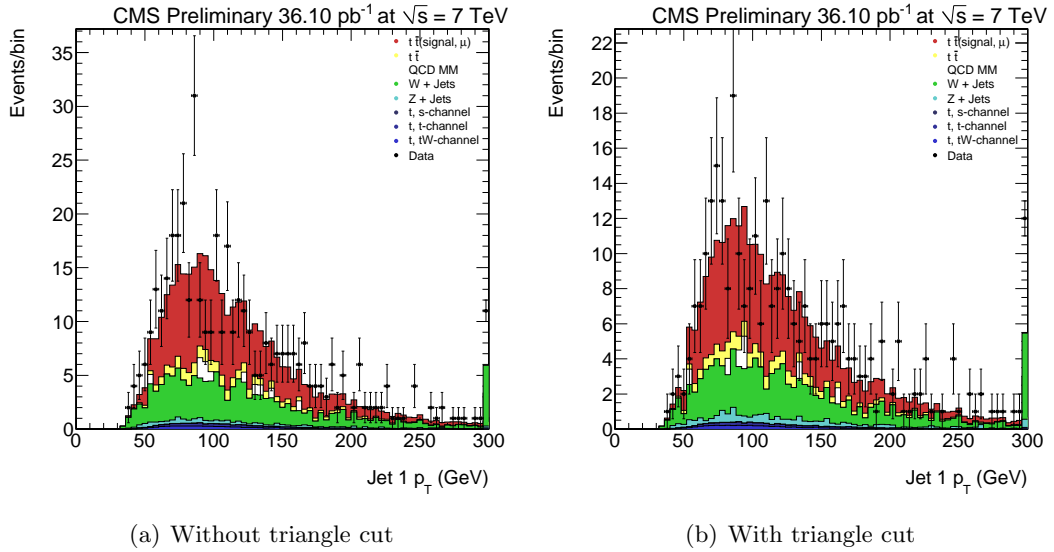


Figure 21: Distribution of the leading jet transverse momentum in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

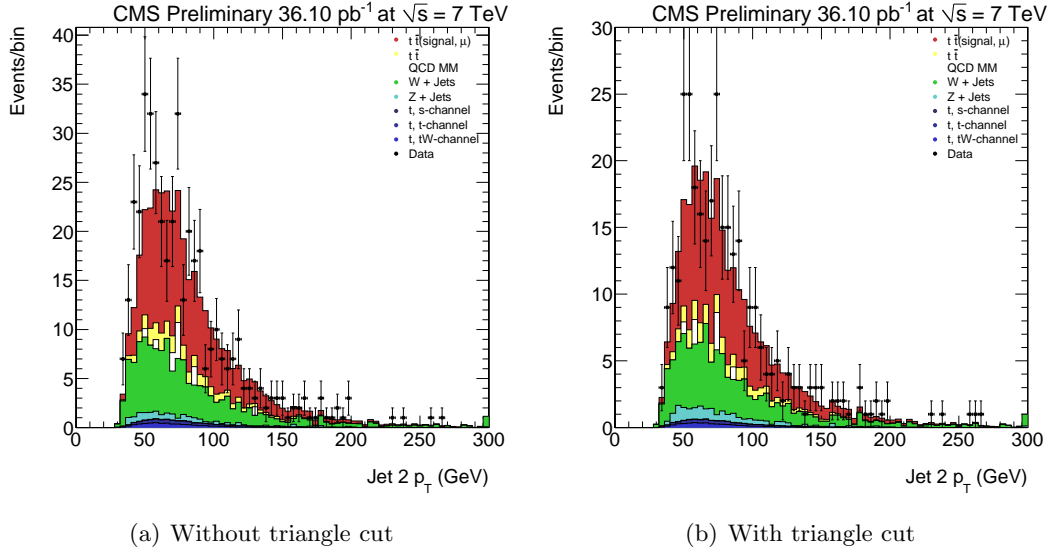


Figure 22: Distribution of the second leading jet transverse momentum in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

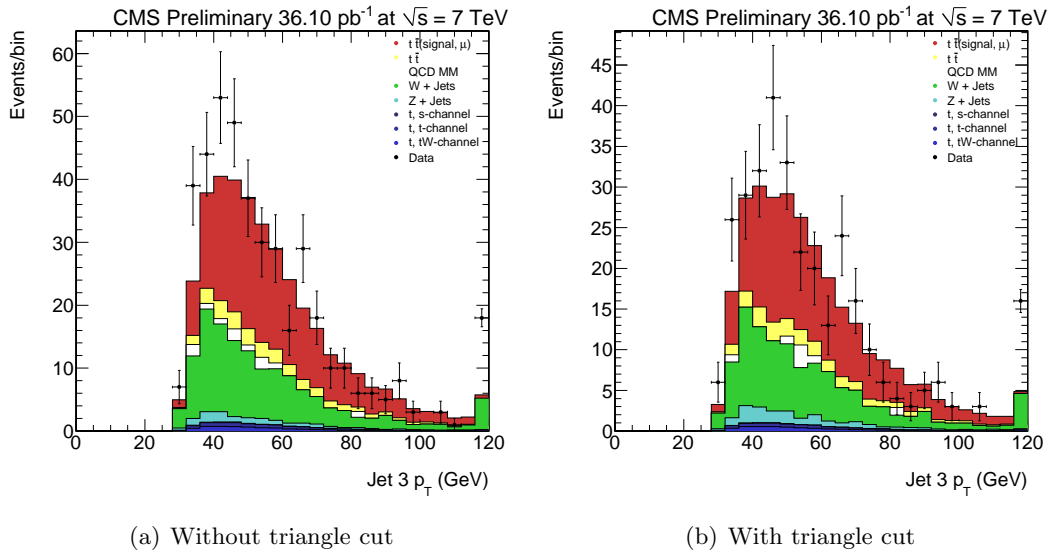


Figure 23: Distribution of the third leading jet transverse momentum in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.

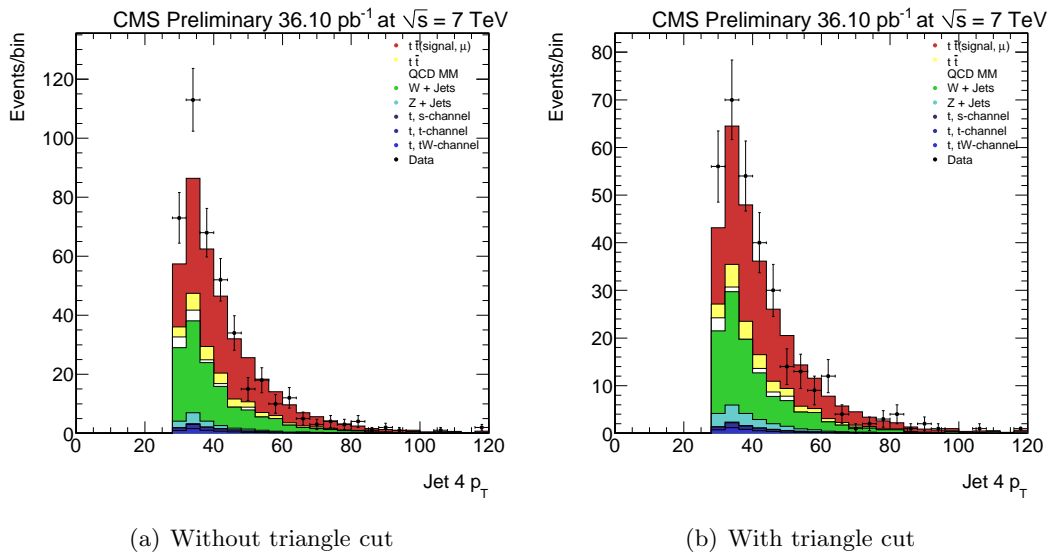


Figure 24: Distribution of the fourth leading jet transverse momentum in the muon channel (a) before applying a triangle cut and (b) after applying the triangle cut as described in the text.