Track/Shower Discrimination for protoDUNE using BDT approach

Mousam Rai 26th August 2020 / ProtoDUNE Meeting/ Supervisor – Dr John Marshall





1



Roadmap

- Terminology
- Small note on vertexing
- BDT variables
- BDT setup
- Brief study on tunable parameters
- Performance plots and numbers
- Summary

Terminology							
		True Co					
	Total PFOs	True Tracks (TT)	True Showers (TS)	/ERSITY OF WARWICK			
Predicted	Predicted Tracks (PT)	True Positive (TP)	False Positive (FP)				
Condition	Predicted Showers (PS)	False Negative (FN)	True Negative (TN)				

- 2x2 confusion matrix
- Sensitivity/True Positive Rate (TPR) = $\frac{TP}{TT}$
- Specificity/True Negative Rate (TNR) = $\frac{TN}{TS}$
- Overall Accuracy = $\frac{TP+TN}{Total PFOs}$
- Trying to maximize Overall Accuracy





- Incoming pion has 2 vertex
 - Interaction Vertex
 - Start vertex
- Use Interaction Vertex as a basis for BDT variables (more on these on the next slide)

BDT Variables

- 13 variables
- Topological Hit spread information
- Calorimetric charge-based information
- Hierarchy hierarchy-based information



Micro	DUNE FD	
Topological	Calorimetric	<u>Hierarchy</u>
length	charge1	nAllDaughter
diff	charge2	nHits3DDaughterTotal
gap		daughterParentNhitsRatio
rms		
diffAngle		
pca1		
pca2		
vertexDistance		

BDT Setup

- Training
 - MCC12 samples WITH space charge prepared by Andy Chappell. (Thanks!)
 - Use tuned BDT parameters (in next few slides)
 - Completeness and Purity Cuts (in next few slides)
 - Fiducial volume cuts : 30 cm in both x-direction, 20 cm in both y-direction, 100 cm in both z-direction
- Testing
 - MCC12 samples WITH space charge prepared by Andy Chappell. (Thanks again!)
 - No completeness or purity cuts
 - Fiducial Volume cuts : see above

Total #Event Files	Train (40% Total)	Test (60% Total)
196	78	118
49	20	29
396	158	238
394	158	236
304	122	182
1339	536	803
229,099	95,229	133,870
86,479	36,154	50,325 6
	Total #Event Files 196 49 396 394 304 1339 229,099 86,479	Total #Event Files Train (40% Total) 196 78 49 20 396 158 394 158 304 122 1339 536 229,099 95,229 86,479 36,154





Brief study on tunable parameters

- We can try to tune certain parameters to increase BDT performance
- Focused on Ntrees, MaxDepth, MinNodeSize, and AdaBoostBeta
- Also investigate completeness and purity cuts used in the training samples



MaxDepth

- Max depth of the decision tree allowed
- In general, the greater the max depth, the better the performance
- But takes longer to train and test

NTrees

- Number of trees in the forest
- As with MaxDepth, in general, more is better
- But also increases training and testing times

MinNodeSize

- Minimum percentage of training events required in a leaf node
- For example, TMVA user guide has a default value of 5%
- But other sources online recommends something like ~1%
- Just need to tune it for your use case

AdaBoostBeta

- Learning rate if using AdaBoost as boosting type
- Based on exponential loss function
- More on CERN TMVA User Guide (Section 7.1)
- Generally, smaller is better but also slower



Tuning Completeness and Purity cuts

- Lesson learnt from working on DUNE FD mc events
- Initially thought training on high completeness and purity PFOs was better but on further study, it turns out, it isn't.
- Mainly because there are a lot of low completeness and purity PFOs and have very sparse hit spread so they don't look anything like high completeness and purity PFOs so what you are training on isn't representative of the actual topologies present in your event.
- Using completeness and purity cuts ONLY during training









Best performance value



Tuned BDT for protoDUNE

 Take all the best performing parameter values and cuts and train a new BDT



	Sensitivity	Specificity	Accuracy	#PFOs
CutFlow	0.4464	0.9524	0.6224	~50,000
Tuned BDT	0.8718	0.7667	0.8352	~50,000

Summary

- ROOT TMVA BDT is significantly outperforming the current cutflow approach (especially at the low nHits region).
- Next steps include implementing the BDT using SKLearn (which is what PANDORA supports), testing and replicating performance numbers and plots, and ideally, work towards its release.
- Any questions and comments are deeply appreciated.



BACK UP SLIDES

Correlation Matrix (signal)





Correlation Matrix (background)

							L	inear co	rrelation	coefficie	ents in %			100
daughterParentNhitsRatio	-1	-2		-1	-1		-1	-1	-1	2	26	64	100	100
nHits3DDaughterTotal	7	-1		6	-4				5	-1		100	64	80
nAllDaughter	11	-1		11	-7		2	3	7	-1	100	65	26	60
charge2	-35	-9		-33	13	-1	-10	-9	-39	100	-1	-1	2	
charge1	51	20	17	52	-16	8	24	23	100	-39	7	5	-1	40
pca2	23	37	10	46	-6	4	51	100	23	-9	3		-1	20
pca1	37	23	1	48	-10	8	100	51	24	-10	2		-1	0
diffAngle	5	10	12	5	-4	100	8	4	8	-1				
vertexDistance	-17		-3	-17	100	-4	-10	-6	-16	13	-7	-4	-1	20
rms	83	19	3	100	-17	5	48	46	52	-33	11	6	-1	-40
gap	7	25	100	3	-3	12	1	10	17		1			-60
diff	9	100	25	19		10	23	37	20	-9	-1	-1	-2	
length	100	9	7	83	-17	5	37	23	51	-35	11	7	-1	10
length diff 9ap rms verlexDistance diffAngle Pca1 Pca2 charge1 charge2 nAllDaughterParentWhiteD						erp _{arentNhitsPar}								



Variable	Definition
length	3D length of a PFO
diff	Average mean difference between the position of hits and a straight line divided by the straight line length
gap	Average max gap distance divided by the straight line length
rms	Average root mean square of linear sliding fit divided by straight line length
vertexDistance	Distance between the PFO vertex and the interaction vertex
diffAngle	Difference between the opening and closing angles calculated over 50% of the pfo closest and furthest from the vertex
pca1	Ratio between the second largest and the largest PCA eigenvalue
pca2	Ratio between the third largest and the largest PCA eigenvalue
charge1	Ratio between sigmaCharge (= (charge – meanCharge)^2) and the mean charge in the collection plane
charge2	Ratio of charge in the last 10% of the pfo and the mean charge in the collection plane
nAllDaughter	Number of all downstream daughter PFOs
nHits3DDaughterTotal	Number of 3D hits in all downstream daughter PFOs
daughterParentNhitsRatio	Ratio between the 3D hits in all downstream PFOs and the parent PFO