Approaches for Assessing Performance of HRMS-based NTA Methods

Christine M. Fisher (O'Donnell)^{1,} <u>Katherine T. Peter²</u>, Seth R. Newton³, Andrew J. Schaub⁴, and Jon R. Sobus³

¹FDA, ²UW Tacoma (current); NIST (former), ³EPA, ⁴SwRI

May 23, 2022 SETAC NTA FTM





NL







Acknowledgements

BP4NTA:

Co-chairs 2020-21 Ben Place (*NIST*) Elin Ulrich (*EPA*)

Co-chairs 2021-22 Christine Fisher (O'Donnell) (FDA) Ruth Marfil-Vega (Shimadzu)

Clearance Reviewers:

Benjamin Place (NIST) James McCord (EPA) Ann Knolhoff (FDA)



Additional members available at: <u>www.nontargetedanalysis.org</u>

(as well as many other great NTA resources!)

The views expressed in this presentation are those of the author(s) and do not necessarily represent the views or the policies of the FDA, NIST, EPA, or SWRI

Why NTA performance assessment?

Targeted analytical methods Performance assessment well established (accreditation!) Familiar terms and metrics to quantify uncertainty – precision, sensitivity, LOD, selectivity...

Readily interpreted, quantitative "final" output

Non-targeted analytical methods Lab-by-lab approaches to QA/QC & performance metrics Lots of uncertainty without clear descriptors [e.g., why was a compound not detected?]

Results often considered "nonfinal" – funneling toward targeted methods

Performance assessment vs. QA/QC



QA/QC happens throughout the NTA workflow...

Performance assessment happens at the END of the workflow!

(Can always do a "look back" to figure out root cause of poor performance)

Let's focus on 3 NTA study objectives



echo targeted analysis approaches?

the confusion matrix (with caveats...)

Fisher et al. Submitted, 2022

Setting the stage: Confusion matrix 101

			Real Condition		
			Honey Classification:		
			Adulterated	Authentic	
Reported Condition	Honey Classification:	Reported Adulterated	True Positive (TP)	False Positive (FP) Type I Error	
			TP = 10	FP = 8	
		Reported Authentic	False Negative (FN) Type II Error	True Negative (TN)	
			FN = 2	TN = 30	



Confusion matrices are routinely applied to assess performance of tests with discrete, often binary, qualitative outputs

Requirements:

- Define positive and negative conditions (positive = rarer)
- A discrete "boundary"
- Sufficient statistical power (# of samples)

Setting the stage: Confusion matrix 101

			Real C	Condition		
			Honey Cl	assification:		
_			Adulterated	Authentic		
Reported Condition	ication:	Reported Adulterated	True Positive (TP)	False Positive (FP) Type I Error	Precision TP TP + FP	False Discovery Rate (FDR)
	ssif		TP = 10	FP = 8	Precision = 0.56	FDR = 0.44
	oney Cla	Reported Authentic	False Negative (FN) Type II Error	True Negative (TN)		
	Ť		FN = 2	TN = 30		
Fisher et al. Submitted, 2022		al. 2d.	True Positive Rate (TPR; Recall, Sensitivity)	False Positive Rate (FPR; Fall-out Rate)	F ₁ Score	Accuracy
		,	TP TP + FN	FP FP + TN	2 × Precision × Precision +	$\frac{\text{TPR}}{\text{TPR}} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$
			TPR = 0.83	FPR = 0.21	F ₁ = 0.67	Accuracy = 0.80
		False Negative Rate (FNR; Miss Rate)		True Negative Rate (TNR; Specificity, Selectivity)	Matthew's Correlation Coefficient (MCC)	
			FN TP + FN	TN FP + TN	$\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}$	
			FNR = 0.17	TNR = 0.79	MCC = 0.55	

LOTS of associated metrics – can be viewed as pairs/groups used interchangeably or in tandem

Note of caution: There is overlapping terminology for targeted analytical methods vs. the confusion matrix!

Assessing <u>sample classification</u> performance with the confusion matrix

Good news: Given a test set of samples, there's a clearly defined CM boundary!



Developing robust sample classification models is challenging – introduces caveats to consider during performance assessment

- Is the test/training set representative? Large enough N? Wellmatched to potential variability? Balanced positive/negative condition?
- Risk of overfitting (1000's of HRMS features, 10's of samples)
- Are feature lists reproducible? Changes over time/across instruments impact model outputs...

Using the confusion matrix for <u>chemical</u> <u>identification</u> performance is trickier...

Need a discrete number of considered chemicals to define the boundary of the confusion matrix

Boundary		Chemical is		
D	n = X	spiked into sample	not spiked into sample	
al is	reported in sample	TP	FP	
Chemic	not reported in sample	FN	TN	

The "chemical universe" is too large...TNs become infinite!

We propose two options for discrete boundaries: Chemicals known to be present and/or reported in a sample A suspect screening database

Boundary 1: Chemicals known to be present/reported in a test sample



- Relies on a spiked test sample
- No TNs! Can't define # chemicals NOT in the sample
- Smaller suite of performance metrics...



Fisher et al. Submitted, 2022

Boundary 2: A suspect screening database



Boundary n = 10,000,000		Chemical is			
		spiked into sample	not spiked into sample		
al is	reported in sample	TP 175	FP 75	Precision 0.70	FDR 0.30
Chemic	not reported in sample	FN 325	TN 9,999,425		
		TPR 0.35	FPR 0.00001	F ₁ 0.47	Accuracy 0.99996
		FNR 0.65	TNR 0.99999	MCC 0.49	

- Relies on a spiked test sample
- Can count TNs! But watch out for huge databases – lots of TNs might not be method-amenable!
- Watch out for metrics that use TN! And note that using F1 Score and MCC alongside Accuracy provides better indication of method performance

Fisher et al. Submitted, 2022

More to consider with ID performance...

• What's actually in the sample??

Rely on <u>initial characterization</u> of the test sample to define TPs.

• What about impurities?? FPs or "unintentional TPs" (uTPs)? Treat them all as FPs to be conservative.

• How to address poor performance?

Start backtracking to evaluate individual workflow steps, use QA/QC to inform where things went wrong...

• What about identification confidence?

Enumerate # of TPs/FPs at each confidence level. Doesn't change metrics, but gives nuance to their interpretation.

 What if there are two potential structures for a given feature? Pick one to report – necessary for confusion matrix. Much room for development....

Defining Quantitative NTA (qNTA)

Generating concentration estimates in the absence of reference standards for the compounds of interest



Chemical Concentration

Fisher et al. Submitted, 2022

The key to assessing qNTA performance: evaluating uncertainty

Easy, but limited predictive value: Calculate error in known test samples, apply to all subsequent results:

[Estimated] [Known]

Harder, but has more predictive value:

- Develop a distribution of response factors for surrogate "calibrants"
 - Incorporate additional info (e.g., predicted ionization efficiency) to narrow which surrogates are used [Groff et al. In review, 2022]



Many remaining challenges, e.g.:

- Accounting for sample recovery & matrix effects?
- Transferability across instruments?

Where do we go from here?

- Current BP4NTA goal is to improve awareness of needs and challenges in NTA performance assessment – not an endpoint, these are not hard & fast recommendations
- What is "good" performance? No one-size-fits-all answer... Different projects and outputs will merit different approaches and stringencies
- Community discussion and further development needed
 - Standardized sample/chemical test sets? Or guidelines?
 - Evaluate transferability of sample classification models
 - Guidance for performance assessment across ID confidence levels
 - Methods to bound qNTA predictions

Thank you!

Kathy Peter ktpeter@uw.edu Other BP4NTA talks/posters to check out:

Ann Knolhoff – *QC chemical mixture* Ruth Marfil-Vega – *Harmonization efforts* Gabby Black – *Chemical space coverage* Jon Sobus – *SRT poster*

Interested in joining BP4NTA or receiving our mailing list?

Head to our website: <u>www.nontargetedanalysis.org</u> Or reach out directly to leadership:

> Christine Fisher (<u>Christine.ODonnell@fda.hhs.gov</u>) Ruth Marfil-Vega (<u>rmmarfilvega@shimadzu.com</u>)

