



www.sbvimprover.com

Strategy for Scoring Prediction Performance

Raquel Norel, PhD

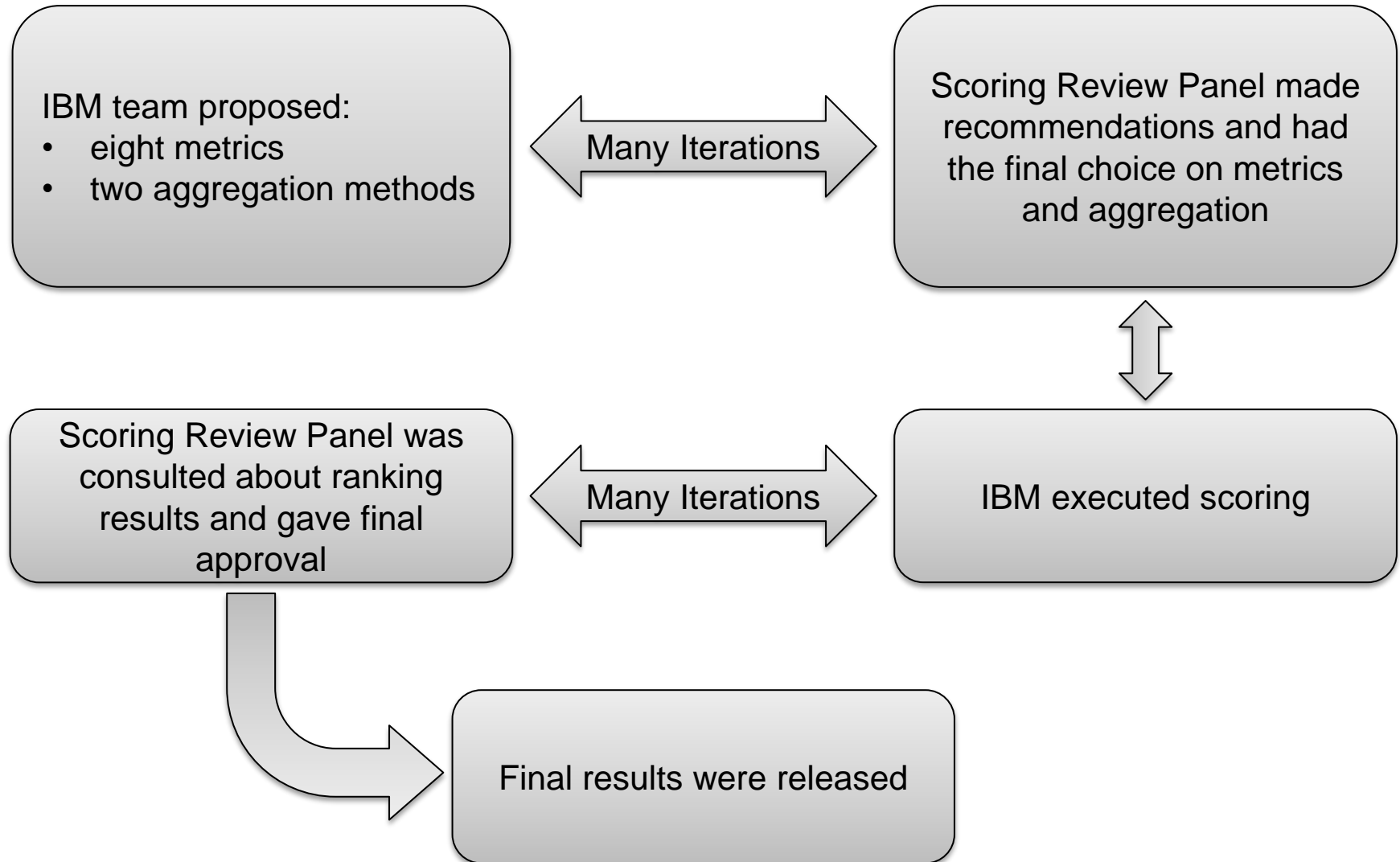
Functional Genomics & Systems
Biology

IBM Computational Biology Center

Outline

- Scoring Process: The Scoring Review Panel
- Criteria for scoring and description of metrics

Scoring Process: at all times the identities of the teams were anonymized

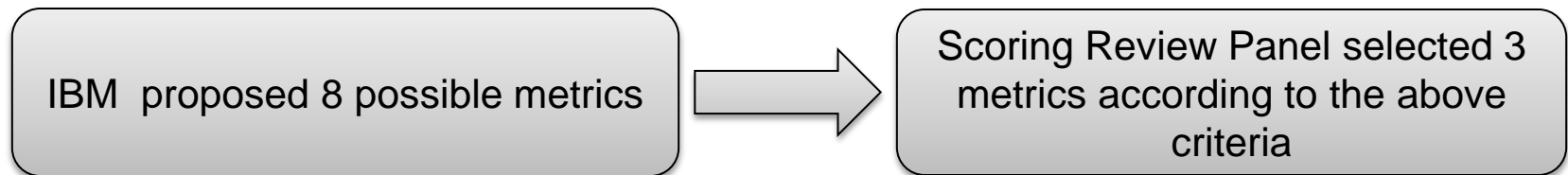


- Scoring Review Panel:
 - **Leonidas Alexopoulos**, National Technical University of Athens / ProtATonce
 - **Jim Costello**, Harvard Medical School
 - **Rudiyanto Gunawan**, Swiss Federal Institute of Technology (ETH) Zurich
 - **Torsten Schwede**, Biozentrum University of Basel & Swiss Institute for Bioinformatics
 - **Alfonso Valencia**, Spanish National Bioinformatics Institute (INB)
- IBM Scoring Team
 - Erhan Bilal
 - Raquel Norel
 - Gustavo Stolovitzky

- Scoring Process: The Scoring Review Panel
- **Criteria for scoring and description of metrics**

Rationale behind the chosen scoring methodology

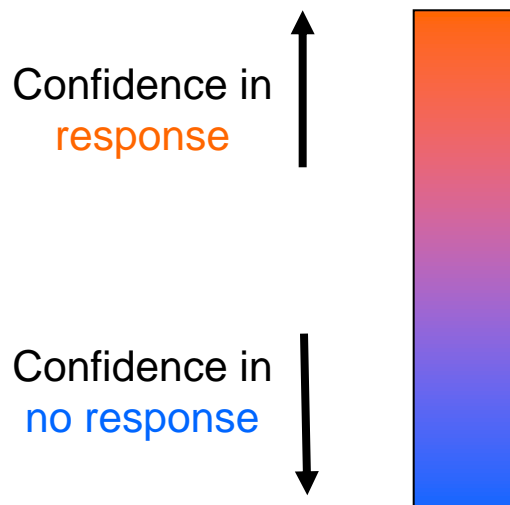
- **Basic premise:** no single metric can capture all the subtleties of a prediction
- We used non-redundant metrics that highlight different qualities of a prediction
 - ✓ Threshold vs. non-threshold
 - ✓ Order-based vs. confidence based
 - ✓ Different ways of rewarding correct vs. incorrect predictions
- A metric should avoid rewarding pathological cases (e.g., predict all items to be one class)
- Unbalanced classes complicated choice of scoring metrics



Binary classification of ordered lists

Participants are required to give confidence values for their predictions of:

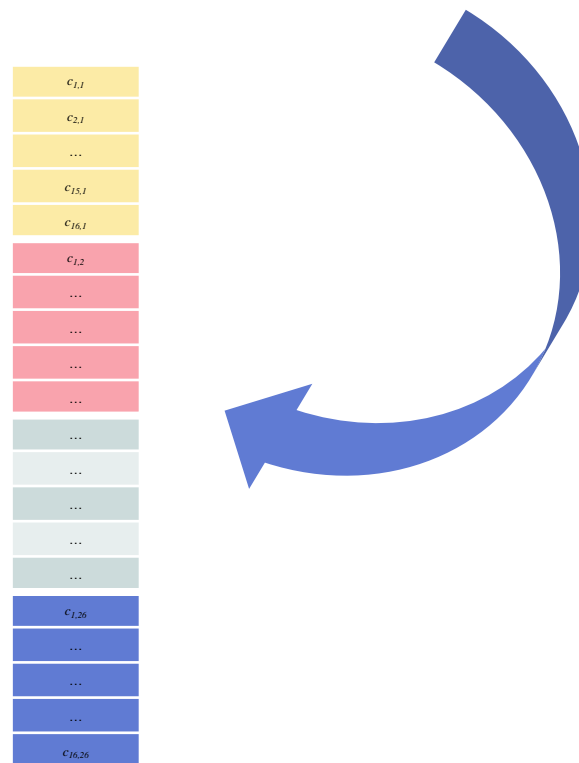
- Protein phosphorylation response to given a stimuli in SC1 and SC2
- Gene set activation for a given stimuli in SC3.



ITEM	Response confidence
item_13	1.00
item_4	0.99
item_52	0.95
...	
item_3	0.01

Vectorize the prediction matrix (few responses or positives calls)

Compound	Cmpd 01	Cmpd 02	Cmpd 26
AKT1	$c_{1,1}$	$c_{1,2}$	$c_{1,26}$
CREB1	$c_{2,1}$
...
TF65	$c_{15,1}$
WNK1	$c_{16,1}$	$c_{16,26}$



Confusion matrix: useful for binary classification

		actual value		total
		p	n	
prediction outcome	p'	True Positive	False Positive	P'
	n'	False Negative	True Negative	N'
total		P	N	

TP \times hit

TN \times correct rejection

FP \times false alarm

FN \times miss

Area Under Precision-Recall Curve (AUPR)

- The precision-recall curve explores, graphically, the tradeoff between these two complementary metrics as k is varied.
- **AUPR** curve presents a single number that summarizes the precision-recall tradeoff.

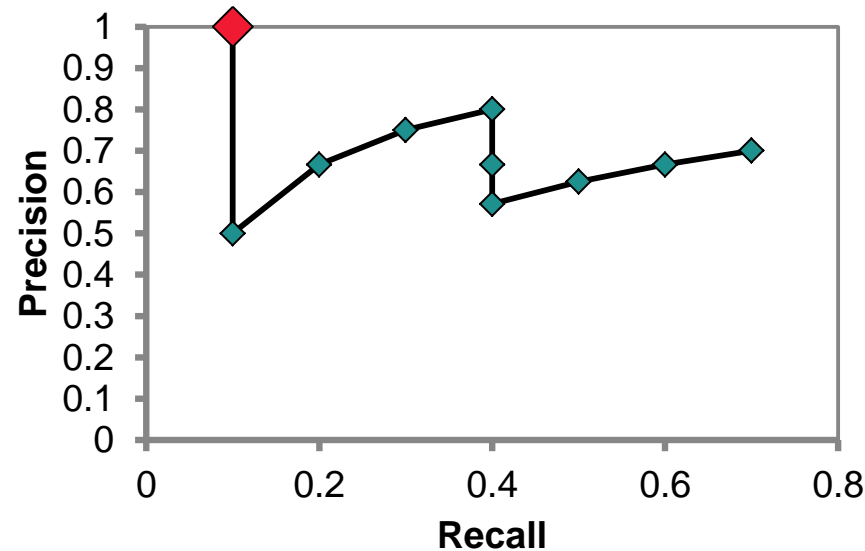
$$prec(k) = \frac{TP(k)}{TP(k) + FP(k)} = \frac{TP(k)}{k}$$

$$rec(k) = \frac{TP(k)}{P}$$

prec: Precision rec: Recall

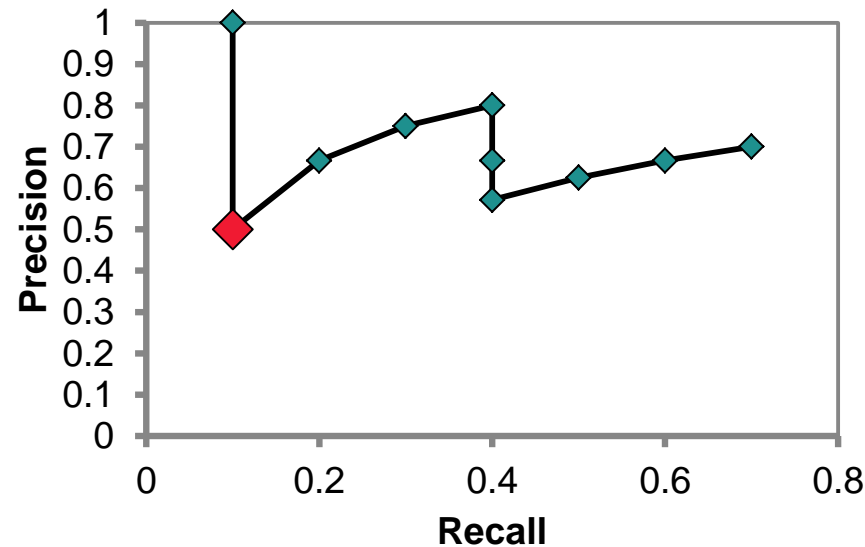
Scoring an ordered list as a binary classification

item ID	Participant predictions	Index k	Gold Standard	Precision TP(k)/k	Recall TP(k)/P
item_4	1.00	1	response	1.00	0.10
item_17	0.99	2	no response	0.50	0.10
item_2	0.95	3	response	0.67	0.20
item_45	0.94	4	response	0.75	0.30
item_13	0.82	5	response	0.80	0.40
...



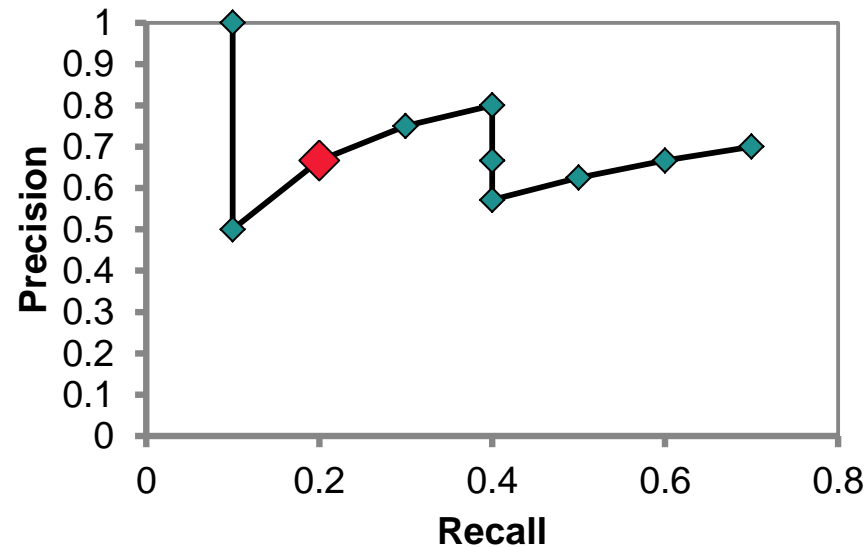
Scoring an ordered list as a binary classification

item ID	Participant predictions	Index k	Gold Standard	Precision TP(k)/k	Recall TP(k)/P
item_4	1.00	1	response	1.00	0.10
item_17	0.99	2	no response	0.50	0.10
item_2	0.95	3	response	0.67	0.20
item_45	0.94	4	response	0.75	0.30
item_13	0.82	5	response	0.80	0.40
...



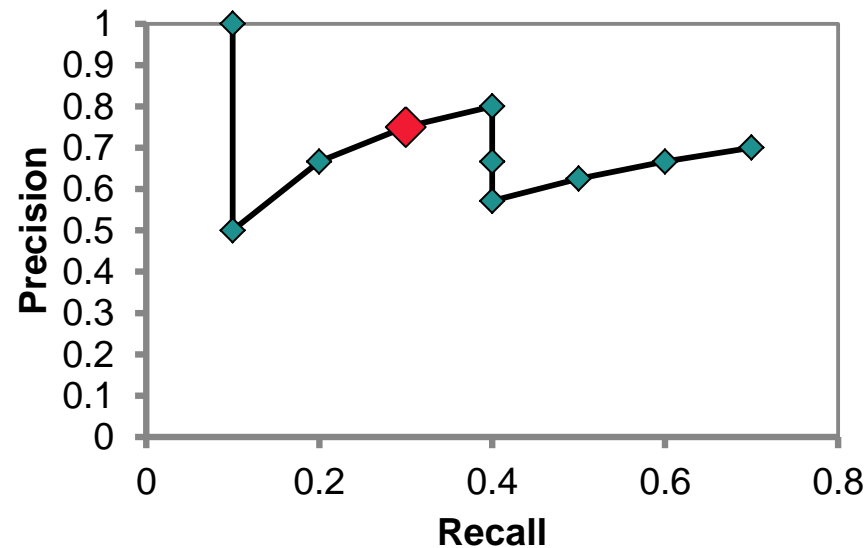
Scoring an ordered list as a binary classification

item ID	Participant predictions	Index k	Gold Standard	Precision TP(k)/k	Recall TP(k)/P
item_4	1.00	1	response	1.00	0.10
item_17	0.99	2	no response	0.50	0.10
item_2	0.95	3	response	0.67	0.20
item_45	0.94	4	response	0.75	0.30
item_13	0.82	5	response	0.80	0.40
...



Scoring an ordered list as a binary classification

item ID	Participant predictions	Index k	Gold Standard	Precision TP(k)/k	Recall TP(k)/P
item_4	1.00	1	response	1.00	0.10
item_17	0.99	2	no response	0.50	0.10
item_2	0.95	3	response	0.67	0.20
item_45	0.94	4	response	0.75	0.30
item_13	0.82	5	response	0.80	0.40
...



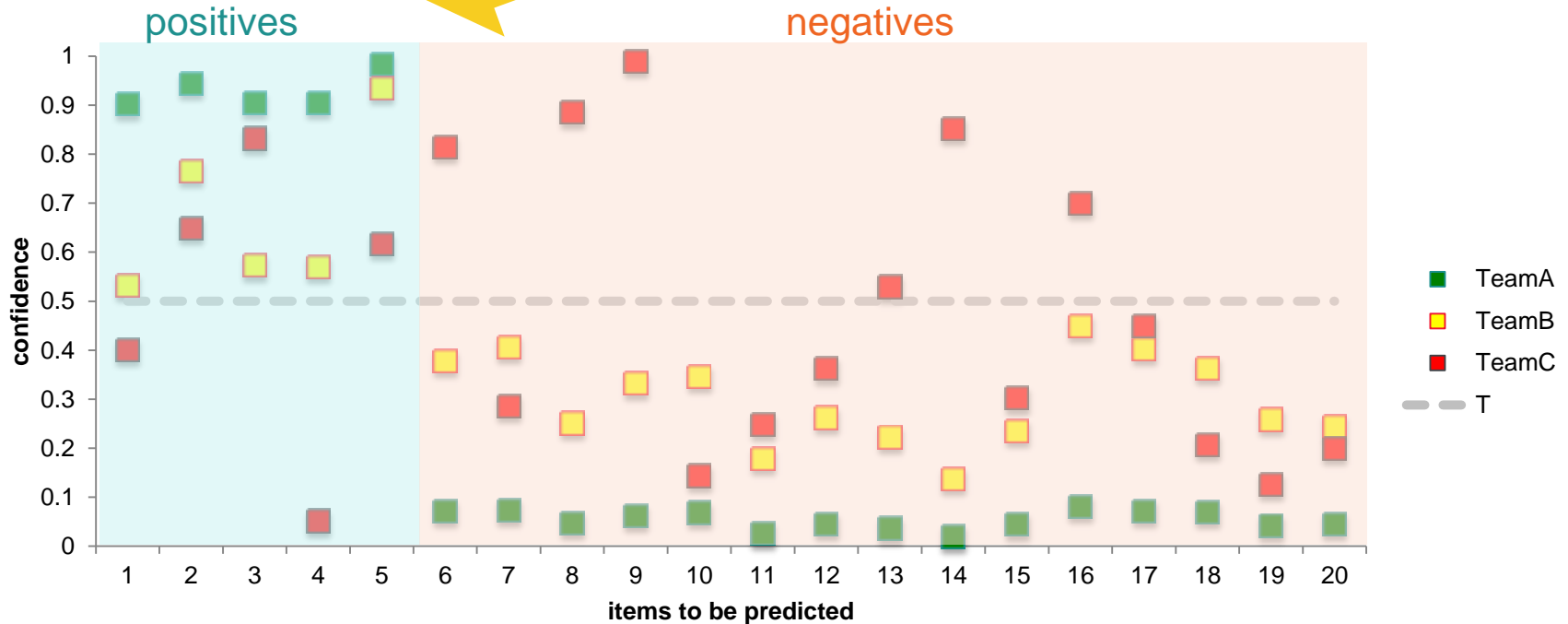
Balanced Accuracy (BAC)

$$BAC = \frac{1}{2} \left(\frac{1}{P} \sum_{C_{ij} \geq 0.5} \frac{1}{|Z_{ij}|^{3.5}} + \frac{1}{N} \sum_{C_{ij} < 0.5} \frac{1}{|Z_{ij}|^{3.5}} \right) = \frac{1}{2} \left(\frac{TP}{P} + \frac{TN}{N} \right)$$

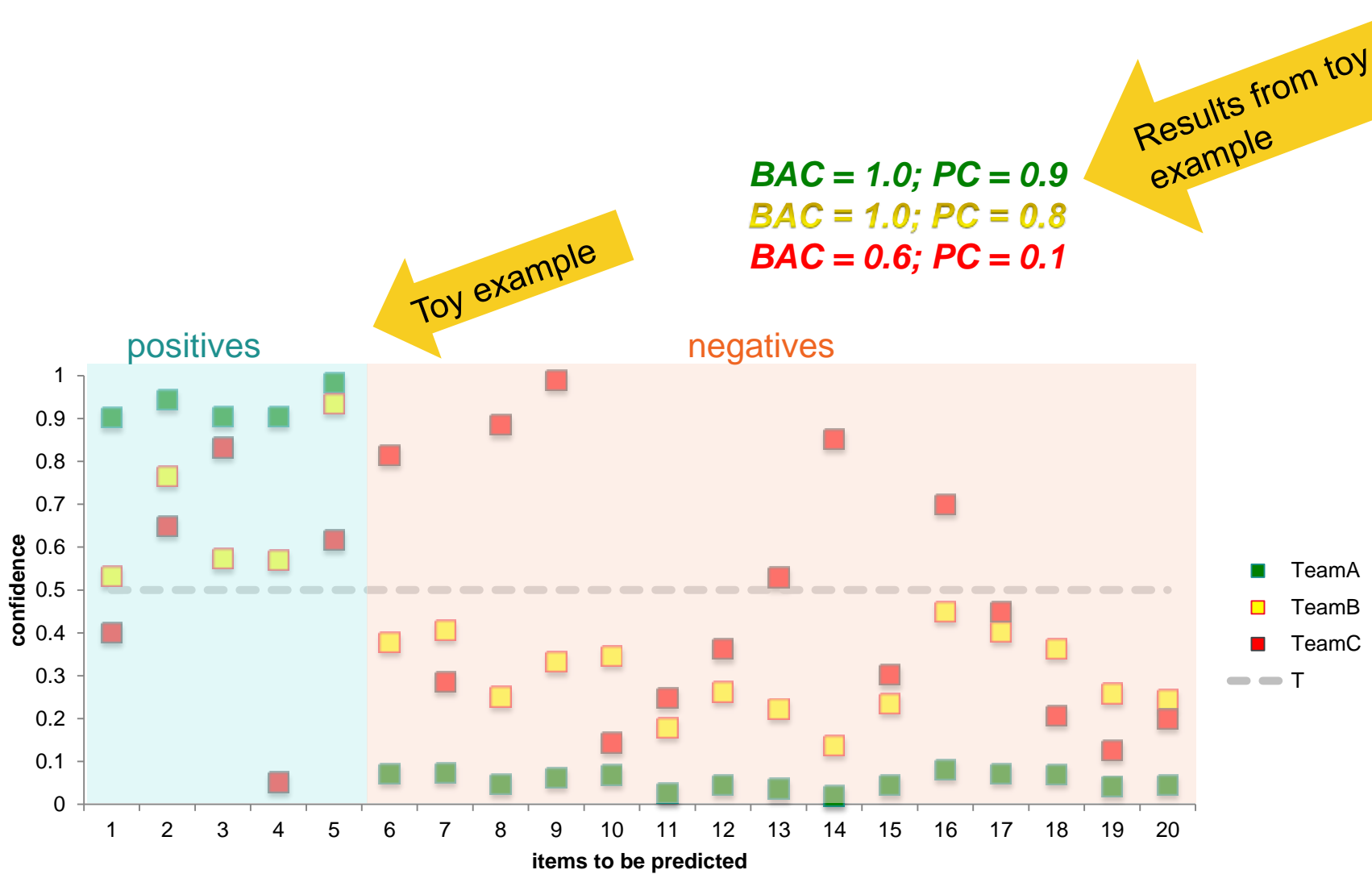
BAC = 1.0
BAC = 1.0
BAC = 0.6

Results from toy example

Toy example



Balanced Accuracy (BAC) + Pearson Correlation (PC)



	AUPR	BAC	Pearson
Threshold	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
Confidence	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
Rank	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

- Random average: $\langle \text{AUPR} \rangle = P/(P+N)$
- Random average: $\langle \text{BAC} \rangle = 0.5$
- Random average: $\langle \text{Pearson} \rangle = 0$
- Random average: $\langle \text{Pearson}_{\text{normalized}} \rangle = 0.5$

To make comparison among metrics easier, we normalized the Pearson correlation:

$$\text{Pearson}_{\text{normalized}} = 0.5 * (\text{Pearson} + 1)$$

With this modification to the Person correlation, all metrics range from [0,1], with 1 being a perfect score.

Criteria for aggregation of performance metrics:

- Final performance must reflect a consensus of best performers in majority of metrics
- All metrics carry the same weight

IBM team proposed two aggregation methods:

- Rank-based
- p -value based



Scoring Review Panel selected:

- Rank-based aggregation across metrics

Conclusions

- Combination of 3 scoring methods
- The scoring methods cover Confidence, Threshold and Rank
- Balanced Accuracy is especially important as it takes into account the imbalance of the sets
- Overall, methods avoid rewarding pathological submissions (“All on”, “All off”)

Thank you for your attention

The sbv IMPROVER project and www.sbvimprover.com are part of a collaboration designed to enable scientists to learn about and contribute to the development of a new crowd sourcing method for verification of scientific data and results. The project team includes scientists from Philip Morris International's (PMI) Research and Development department and IBM's Thomas J. Watson Research Center. The project is funded by PMI.