SYSTEMS BIOLOGY VERIFICATION

www.sbvimprover.com

The Systems Toxicology Computational Challenge:

Marker of Exposure Response Identification

Scoring results & Lessons learned

Vincenzo Belcastro, PhD July 11th 2016, ISMB, Orlando



- Participation & Scoring methodology
- Scoring Results & Top performing teams
 - Sub-Challenge 1 (SC1)
 - Sub-Challenge 2 (SC2)
- Post-challenge Analysis Initial results
 - Misclassifications
 - Confidence values distribution per class
 - Gene signatures co-occurrence for smoking exposure response status (Smokers vs Non-Current Smokers)
 - Human-specific signatures
 - Species-independent signatures

Conclusions



© 21



Participation & Scoring methodology



Computational Challenge Participation across the globe





135 registered participants in 61 teams - 23 teams submitted files in at least 1 sub-challenge

Datasets provided for training and tests

Gene expression data generated from human and mouse blood samples





S/3R4F: Smokers / 3R4F (exposure to smoke from a reference cigarette) FS/Cess: Former smokers / Cessation NS/Sham: Never smokers / Sham NCS: Non-current smoker

Freedom to use two separate models for 2-class prediction for each step, or directly a 3-class prediction model

https://sbvimprover.com/challenge-4/the-computational-challenge/hbs

Scoring metrics and aggregation





Qualifications rules & procedure for scoring



Rules as per challenge description

- All requested files submitted:
 - o 4 prediction files
 - o 2 gene lists (in case of 3-class model: 1 gene list)
 - o 1 write-up
- Confidence values for common samples are the same
- Gene list length does not exceed 40 genes

Scoring procedure

- Anonymized submissions
- Results and final ranking presented to and approved by a external and independent Scoring Review Panel





Scoring Results Sub-challenge 1

© 2^r PROVE

Participants' prediction score ranking



23 teams provided submissions \rightarrow 12 teams had valid submissions



And the winners are...







TEST TRAIN dset2 dset5 dset1 dset4 **3R4F** S 3R4F S (109) (40) (27) (12) **Scoring Results** FS Cess FS Cess Sub-challenge 2 (57) (27) (26) (8) Sham NS NS Sham (45) (58) (28) (13)

© 20 PROVE

Participants' prediction scoring



15 teams provided submissions \rightarrow 6 teams with valid submissions



© 2016 sbv IMPROVER

And the winners are...





• Team 219

- Omer Sinan Sarac
- Ismail Bilgen
- Ali Tugrul Balci
- Team 250
 - Rahul Kumar
 - Sandeep Dhanda
- Team 264
 - Adi L Tarca
 - Roberto Romero







USA



Post-challenge analysis

Misclassification sub-challenge 1



Misclassification Sub-Challenge 1



Smokers vs Non-Current Smokers

- **Best performers have perfect predictions** for smoking exposure Samples from former smokers tend to be more frequently misclassified 269 215 250 290 283 222 257 247 221
- Correct
 Misclassified
- Smoker
- Former Smoker
- Never Smokers

Misclassification – Sub-Challenge 1



Former Smokers vs Never Smokers

- Very difficult to prediction cessation status 222 283 215 257 247 269 290 250 221 © 2016 sby IMPROVER
- Correct
 Misclassi
- Misclassified
- Former Smoker
- Never Smoker



Post-challenge analysis

Misclassification sub-challenge 2



17

Misclassification – Sub-Challenge 2



Smokers vs Non-Current Smokers



- Correct
- Misclassified
- Smoker/3R4F
- Former Smoker/Cessation
- Never Smokers/Sham

Misclassification – Sub-Challenge 2



Former Smokers vs Never Smokers



- Correct
- Misclassified
- Former Smoker/Cessation
- Never Smokers/Sham



Post-challenge analysis

Confidence values distributions per class



20

Confidence value distributions per class for all and top 3 teams Sub-challenge 1





<u>Test dataset</u>: human data from an independent study



© 2016 sbv IMPROVER

Confidence value distributions per class for all and top 3 teams Sub-challenge 2







Post-challenge analysis

Gene signatures co-occurrence for smoking exposure response (Smokers vs Non-Current Smokers)

© 20 PROVE

Human-specific gene signature for smoke exposure response Sby MPROVER Sub-Challenge 1



* Matching based on gene name

Conclusions



- Successful **worldwide participation** to the challenge
- Gene expression changes measured in **blood are informative of exposure status**:
 - 1. Identification of **smoking exposure** status is **possible** (S vs NCS)
 - 2. Identification of **cessation status** is more **challenging** (FS vs NS)
- Participants succeeded in development of **inductive classification models**
- Random forest, linear discriminant analysis, partial least square discriminant analysis and logistic regression were machine learning methods used by best performing teams
- Samples from former smokers tend to be more frequently misclassified
- Exposure response markers predictive of smoking status were **identified in human blood** and included a core gene subset **highly consistent** across teams



sbv IMPROVER team

- Vincenzo Belcastro
- Filipe Bonjour
- Stephanie Boue
- Laure Cannesson
- Anouk Ertan
- Sylvain Gubian
- Julia Hoeng
- Florian Martin
- David Page
- Manuel Peitsch
- Carine Poussin
- Alain Sewer
- Marja Talikka
- Bjoern Titz
- Yang Xiang

To learn more, visit Questions? Contact Us sbvimprover.com/comp-start sbvimprover.RD@pmi.com

The sbv IMPROVER project, the websites and the Symposia are part of a collaborative project 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 is led and funded by Philip Morris International. The current challenges, website and biological network models were developed and are maintained as part of a collaboration with Selventa, Douglas Connect, SBX-Garuda, Nebion, OrangeBus and ADS. For more information on the focus of Philip Morris International's research, please visit www.pmiscience.com.

Scoring Review Panel

- Prof. Leonidas Alexopoulos (National Technical University of Athens)
- Prof. Rudiyanto Gunawan (ETH Zurich)
- Dr. Alberto de la Fuente (Leibniz Institute for Farm Animal Biology)