Applications of Computational Biology and Machine Learning to Immuno-Oncology

Non-invasive approaches for cancer detection and monitoring

Robert B Scharpf

Johns Hopkins University

Department of Oncology

SITC Cancer Immunotherapy Winter School Austin, TX | January 26, 2022



R.B.S. is a co-founder and consultant of Delfi Diagnostics, and owns Delfi Diagnostics stock subject to certain restrictions under university policy. Johns Hopkins University owns equity in Delfi Diagnostics.

Liquid biopsies

- most cell-free DNA (cfDNA) is shed from hematopoietic cells
- for cancer patients, some cfDNA may be shed by tumor cells (ctDNA)



Liquid biopsy approaches



- Leary et al, *Science TM*, 2010
- Leary et al, *Science TM*, 2012
- Phallen et al, Science TM, 2017
- Cristiano et al, Nature, 2019
- Mathios et al, Nature Comm, 2021

Detection of tumor-specific alterations

Continuum of cancer care

Early cancer detection and screening



Screening & Early Detection

Challenges

- de novo detection required
- ctDNA/cfDNA may be very small while disease is still localized
- cancer prevalences in most screening populations are low; positive predictive value of a positive liquid biopsy test may be small

Fraction of ctDNA with mutations is small





Cristiano et al., *Nature*, 2019

cfDNA fragments orginate from nucleosomes



Altered fragment lengths of ctDNA



• Intra-patient comparison of fragment lengths

Negative control



• 20 loci where germline alterations were detected

Genome-wide alterations of fragmentation patterns



A/B compartments (1 - 2x coverage)



Genomic and chromatin changes







Stephen Cristiano^{1,2,15}, Alessandro Leal^{1,15}, Jillian Phallen^{1,15}, Jacob Fiksel^{1,2,15}, Vilmos Adleff¹, Daniel C. Bruhm¹, Sarah Østrup Jensen³, Jamie E. Medina¹, Carolyn Hruban¹, James R. White¹, Doreen N. Palsgrove¹, Noushin Niknafs¹, Valsamo Anagnostou¹, Patrick Forde¹, Jarushka Naidoo¹, Kristen Marrone¹, Julie Brahmer¹, Brian D. Woodward⁴, Hatim Husain⁴, Karlijn L. van Rooijen⁵, Mai-Britt Worm Ørntoft³, Anders Husted Madsen⁶, Cornelis J. H. van de Velde⁷, Marcel Verheij⁸, Annemieke Cats⁹, Cornelis J. A. Punt¹⁰, Geraldine R. Vink⁵, Nicole C. T. van Grieken¹¹, Miriam Koopman⁵, Remond J. A. Fijneman¹², Julia S. Johansen¹³, Hans Jørgen Nielsen¹⁴, Gerrit A. Meijer¹², Claus Lindbjerg Andersen³, Robert B. Scharpf^{1,2*} & Victor E. Velculescu^{1*}

DELFI



Noninvasive cancer screening (DELFI)



High sensitivity and specificity overall



• gradient boosted model





By histology



Statistical challenges and rationale

Genome-wide cfDNA fragmentation data

${\bf Participants}\,\,(n):$

• ~400 (~200 cancer cases and ~200 non-cancers)

• Non-cancers were part of a screening population

• Multiple cancer types: ovarian, lung, breast, and 5 others

cfDNA features (p) :

- ~1000 statistical summaries of cfDNA fragmention
- 39 measures of chromosome arm aneuploidy (\$z\$-scores)
- mitochondrial representation (single number)

Our primary goal is prediction

- Can we distinguish patients with cancer from individuals without cancer?
- We are willing to trade some interpretability for more sensitive and specific detection of cancer
- *n* << *p*

Prediction task is complex

- How do fragmentation profiles across the genome differ from fragmentation profiles we observe in non-cancer populations?
- The combination of features useful for prediction may not necessarily explain important biological pathways or provide mechanistic insight

Machine learning

In machine learning, a model learns from examples rather than being programmed with rules

Rajkomar et al., NEJM, 2019

- Features: the inputs (fragmentation characteristics)
- Labels: cancer or non-cancer

Tradeoff of model complexity and interpretability

n << p:

- "deep learning" approaches will not outperform simpler models
- hypothesis-driven approaches that leverage known biology and features can lead to more useful and replicable models



Mechanics

1. specify resampling-based approach for cross-validation

- 2. specify a machine learning architecture (logistic regression, random forest, etc)
- 3. train and test
- 4. summarize predictive performance
- 5. specify a *final* model

Cross-validation involves a lot of repetition

• Easy to make mistakes

Use existing infrastructure for:

- resampling (**rsample**)
- model specification (tidymodels, caret)
- tuning hyperparameters and nested cross-validation
- training and testing a model (tidymodels, recipes, workflows, and others)
- summarizing performance (tidymodels)

Resampling approach: 5-fold cross validation

```
library(tidyverse)
library(magrittr)
library(rsample)
ngenes <- 200
nsamples <- 50
x <- matrix(rnorm(ngenes*nsamples), ngenes, nsamples)
y <- factor(rbinom(nsamples, 1, 0.5))
dat <- t(x) %>%
    set_colnames(paste0("feature_", seq_len(ncol(.)))) %>%
    as_tibble() %>%
    mutate(y=y)
training.test <- vfold_cv(dat, v=5)</pre>
```



Train and test

Summarize performance

collect_metrics(lr_fit)

##	#	A tibble:	: 2 × 6				
##		.metric	<pre>.estimator</pre>	mean	n	std_err	.config
##		<chr></chr>	<chr></chr>	<dbl></dbl>	<int></int>	<dbl></dbl>	<chr></chr>
##	1	accuracy	binary	0.4	5	0.0447	<pre>Preprocessor1_Model1</pre>
##	2	roc_auc	binary	0.5	5	0	Preprocessor1_Model1

Repeated cross-validation

- Our initial randomization of patients to the 5 folds was arbitrary
 - performance assessment could be slightly biased because of the randomization
- Repeat the process r times and average the predictions across the r repeats

Cross-validation or split-study validation?

- Resampling based methods such as bootstrap or repeated k-fold cross-validation are nearly always your best option

Exceptions:

- You have a large study (20,000+) and it does not bother you that you will not learn anything from 1/3 of the dataset
- Required for regulatory purposes
- Your dataset consists of multiple independent studies
 - consider multi-study models
 - approaches for external validation

Early detection of lung cancer

LUCAS study

Altered fragmentation profiles



Machine learning model

• Machine learner: penalized logistic regression

Features:

- PCA of the fragmentation profiles
- z-scores for chromosomal aneuploidy

278 x 504 matrix



36/49



Dimensionality reduction

- Principal components that explain 90% of the variance across samples
 - 278 x 10 matrix
- Regression coefficients for chromosome arm-level copy number

Stochasticity of model across training sets

- Columns display 50 training sets from repeated crossvalidation
- Rows are regression coefficients



Scores by stage and histology



Internally cross-validated performance



Internally cross-validated performance



External validation

 Fixed model and cutoff at 80% specificity

Sequential screening



Sequential screening



Fragmentation profiles differ between cancer types



Tissue of origin



Disease monitoring



Machine learning for noninvasive detection and monitoring of cancer

- Genome-wide cfDNA fragmentation profiles reflect abnormal packaging and genomic content of cancer genomes
- Liquid biopsies can detect cancer recurrence early and provide opportunities for intervention
- With enough examples, gradient boosted models / random forests can do a reasonable job at feature selection and prediction
- Leverage known biology when possible especially for n << p
 - fragments derived from nucleosomes, aneuploidy
 - coverage at transcription start sites
- Multi-study models, ensembling machine learners, and approaches for external crossvalidation can lead to more replicable classifiers

Acknowledgements

- Vilmos Adleff
- Akshaya Annapragada
- Jacob Carey
- Stephen Cristiano
- Jacob Fiksel

- Alessandro Leal
- Dimitrios Mathios
- Noushin Niknafs
- Jillian Phallen
- Victor Velculescu