

Clifton (Dave) Fuller, MD, PhD Assistant Professor Head & Neck Section

**MDAnderson** Cancer Center

Machine Learning, "Big Data", Informatics, Radiomics, and Clinical Outcomes

The Future of Innovation in Clinical **Radiation Oncology?** 

# C.D. Fuller Acknowledgment/Disclosures

- 2016-17 Funders: The Andrew Sabin Family Fellowship Program, through an endowment established by the Andrew Sabin Family Foundation
- .
- Andrew Sabin Family Foundation <sup>114</sup> Ingram Integrit and Interdent Exactinet of yith A direct gift from the Beach Family of Phoenix, AZ. NIH Big Data to Knowledge (BD2K) Program of the National Cancer Institute (NCI) Early Stage Development of Technologies in Biomedical Computing, Informatics, and Big Data Science Award (1R01CA214825-01) National Science Foundation, Division of Mathematical Sciences, Quantitative Approaches to Biomedical Big Data (QuBBD)/Big Data to Knowledge (BD2K) Program (NSF1557559; National Cancer Institute Early Stage Development of Technologies in Biomedical Computing, Informatics, and Big Data Science (1 R01 CA214825-01) National Cancer Institute and Craniofacial Research (NR56/R01 DE025248-01) National Cancer Institute Grant MD Anderson Head and Neck Specialized Programs of Researcl Excellence (SPORE) Development Award (P50CA937007-10) Elekta ABMD Anderson MR1-LinAc Consortium Seed Grant\*

- Elekta AB Travel support & Honoraria

#### #-Philantropic individuals/agencies/societies

\*-Corporate/industry funders +-Federal or state funding agencies







# Institutional/Departmental Team





# MD Anderson Multi-disciplinary Symptom Working Group





Kate Hutcheson MD, PhD Head and Neck Surgery PhD Speech Pathology



bdallah Mohamed Jihong Wang Dave Fuller MD, MSc PhD MD, PhD Radiation Oncology Radiation Oncology Radiation Oncology

National Institute of Dental and Craniofacial Research (NR56/R01 DE025248-01; SY Lai, PI)

•National Cancer Institute Grant MD Anderson Head and Neck Specialized Programs of Research Excellence (SPORE) Development Award (P50CA097007-10, J Myers, PI)

•National Institutes of Health/National Cancer Institute Grant (R03 CA188162-01A1; KA Hutcheson, PI)

Jinzhong Yang PhD Radiation Oncology











# A few definitions...

#### Informatics [ The Discipline]

 "Informatics is the science of processing data for storage and retrieval; information science as a field."

#### "Big Data" [The Input]

- "Big Data is high-volume, high-velocity and/or high-variety informatic assets that demand cost-effective, innovative forms of information processing that enable enhanced insight, decision making, and proc automation."- Gartner
- Machine learning ("statistical learning") [The Methodology]
- "Machine learning explores the study and construction of algorithms that ca learn from and make predictions on data – such algorithms overcome follov strictly static program instructions by making data-driven predictions or decisions, through building a model from sample inputs."

# "Big Data"

- "Big data is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it too...and most don't do it very well."
  - Dan Ariely







#### Types of Machine Learning - At a glance















Kagadis et al.: Medical physicists and health care applications of informatics











Ira Kalet, PhD



Retired CSE adjunct professor Ira Kalet passed away last night after a long battle with cancer.

Ira joined the University of Washington in 1978 in the then newly formed Department of Radiation Oncology. Subsequently he held adjunct appointments in Computer Science & Engineering, Bioengineering, and Biological Structure, and a joint appointment in Medical Education (now the Department of Biomedical Informatics and Medical Education)

#### SPECIAL ARTICLE

# **Technology for Innovation in Radiation Oncology**

Indrin J. Chetty, PhD,\* Mary K. Martel, PhD,<sup>†</sup> David A. Jaffray, PhD,<sup>‡</sup>

- Integrating radiation oncology databases across the discipline will facilitate science and elevate the quality of care (45). The reation of a Virtuel Chineal Trials Group that enables federated databases at different institutions for conducting cooperative research is a consideration. Sharing practices and outcomes will permit high mean and tight variance in clinical practice and will improve quality (46).
- 2. Tools need to be created and made available for patients Tools need to be created and made available for patients and physicians to discuss treatment routions, as recom-mended by the Patient-Centered Outcome Research Institution. Such an approach will divise the development of metarteatment planning systems, in which one pre-senties an outcome, not a treatment (eg specification of a 95% local control rate at 5 years with 5% grade 3 or more dysprage) (c, 47). This could also be expanded beyond radiation oncology.
- 3. Expertise in the informatics domain among radiation oncology professionals needs to be developed (6). The most suitable candidates with the appropriate skill sets and multidisciplinary knowledge to succeed in this space are likely medical physicists or physicians with strong

computational backgrounds. Training grants for devel-oping programs for oncology informatics will provide these individuals with the knowledge needed to support informatics research initiatives.

A. Informatics toolst manares. A. Informatics toolst mend to be developed to support the monitoring of the quality of oncology care at the point(s) of delivery (4). Real word—based evidence approaches are emerging in other domains and will also benefit the field of radiation oncology. The other-quoted statements that 3% differences in dose result in significant changes in tumor control and normal issue complication proba-hic time of the state of the state of the statements in tumor control and normal issue complication proba-hic time of the statement of the statement of the statements in tumor control and normal issue complication proba-hic time of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement is the statement of the statement of the statement of the statement of the statement is the statement of the st



NATIONAL CANCER INSTITUTE Informatics Technology for Cancer Research

#### **Funding Opportunities**

ITCR has issued four Funding Opportunity Announcements aimed at successive stages of informatics technology development.

#### Algorithm Development

PAR-15-334 Development of Innovative Informatics Methods and Algorithms for Cancer Research and Management (R21)

#### Prototyping & Hardening

PAR-15-332 Early-Stage Development of Informatics Technologies for Cancer Research and Management (U01)

# Enhancement & Dissemination

PAR-15-331 Advanced Development of Informatics Technologies for Cancer Research and Management (U24)

#### -ustumment

PAR-15-333 Sustained Support for Informatics Resources for Cancer Research and Management (U24)

#### Clinical Informatics Becomes a Board-certified Medical Subspecialty Following ABMS Vote

Thursday, September 22, 2011 AMIA to offer prep courses for clinicians who sit for Board Exam

Washington, DC—Today, AMIA—the association for informatics professionals announces the success of a multi-year initiative to elverate clinical informatics to an American Board of Medical Specialities (AMIAS) subspacelity certified by an examination administered by the American Board of Preventive Medicine and available to physicians who have primary speciality certification through the American Board of Medical Specialities. Joining such subspecialities as pediatric



# Fellowship in Clinical Informatics: Radiation Oncology Track

#### Fellowship in Clinical Informatics: Radiation Oncology Track

Clinical informatics is the subspecialty of all medical specialties that transforms health care by analyzing, designing, implementing, and evaluating information and communication systems to improve patient care, enhance access to care, advance individual and population health outcomes, and strengthen the clinicianpatient relationship.



#### School of Medicine Radiation Medicine

# Health Informatics via Machine Learning for the Clinical Management of Patients

D. A. Clifton, K. E. Niehaus, P. Charlton, G. W. Colopy

Conclusions

4

#### Yearb Med Inform 2015;10:38-43 http://dx.doi.org/10.15265/1Y-2015-014 Published online August 13, 2015

We conclude by emphasising that the field of health informatics systems based on machine learning, drawing on disparate datatypes from the ICU, the wider hospital, and from (potentially very complex) EHR data, is in its infancy. While the majority of hospitals in the developed world have implemented EHR systems of some kind, the integrated use of the large quantities of data that arise from such systems is not employed at scale. This



# The curse of clinical implementation: From n=895 to n=5

		TABLE 3   Articles with further processing strategies and approach collected data.		ther processing strategies and approaches of	
		Reference	Year	Summary	
		Brown et al. (16)	2007	Analysis tools connected to data management system for quantitative image analysis in metastatic lung cancer patients; automatic noo detection and segmentation for CAD evaluation detection and segmentation for CAD evaluation	
TABLE 2   Topics of articles.		Carey et al. (17)	2012	Analysis tools used on imaging files stored in	
Торіс	No. of articles			database in lung cancer patients; manual image analysis; no communication standardization mentioned	
System use For clinical trial or analysis For clinical routine Dates incomparison	469 370 99	Haak et al. (11, 16)	2014	Analysis tools connected to EDC system for automatic image and biosignal analysis; communication standards used: web services, COM, SOAR SETP. HTTP	
System comparisons with paper-based standard or other system: System review, recommendations, and issues Not assigned	200 3 24 95 39	Kessel et al. (12, 19)	2012	Analysis tools connected to documentation database via SQL interface; semiautomatic CT image registration and segmentation of pancreatil cancer patients, as well as dose calculation of	
	N = 895			radiation plans; communication standards used: kil 7, DICOM, https://	
		Ozyurt et al. (20)	2010	Analysis tools used on local copies of neuroimagin data after query and download from the data management system; results are transferred back via web services; communication standards used; web services; SDAP; DICOM; https:	
		CAD: computer-aidec medicine; EDC, electi language; CDM, obje	f diagnos tonic data ct data m	is; DICOM, sligital imaging and communications in capture: HL7, heath level 7; SQL, structured query ode: SQAP simple object access protocol.	

# Kessel and Combs Critique

It is not the lack of technology or tools that keep "the health revolution" from coming, but the lack of expertise, specifications, and concepts (4). One of the most common weaknesses found is the lack of standardization. Most researchers create an individual in-house solution without considering communication standards, such as DICOM, HL7, https, and html (37). These solutions work only in their own environment and are tailored to meet their requirements. This might be necessary up to a certain level, as already stated that there is no "one-size-fits-all" solution for documentation of clinical trials, research data, or patient data





# The "ontology" problem Without common terminology, content is obscured...and we may not be aware of it! What we say to dogs Okay Graper Twe had It? Stay out of the garbage or elser What we say to dogs Okay Graper Twe had It? Okay Graper Twe had It?

# Standards and ontologies

- In informatics, a standard is voluntary consensus developed or adopted by voluntary bodies
   Can also include
  - definitions
  - · a classification system or components,
  - · a series of procedures to follow,
  - specifications for sizes or dimensions, materials, performance, design or operations.
    quality measures or describe amounts of materials.
  - process, a system, a service, a practice, or a specific product.
- In computer science and information science, an ontology is a formal naming and definition of the types, properties, and interrelationships of the entities that really or fundamentally exist for a particular domain of discourse.
- Example: DICOM-RT is a standard, TG-263 OAR names are an ontology.
- · Functionally speaking, a relational standard=an ontology.

#### EXAMPLE: Can we just name OARs the same?

#### Let's start by trying to fix the standardization problems for DVH data

- TG 263 Standardizing Nomenclature for Radiation Therapy
- group of 57 stake holders
- domestic and international groups
- representing a broad range of perspectives

Roles	Professional Societies	Clinic Types	Specialty Groups
Physician	ASTRO	Academic	IHE-RO
Physicist	AAPM	Community	Dicom Working Group
Vendor	ESTRO	Large Practice	NRG
Dosimetry		Small Practice	IROC

Slide courtesy of Chuck Mayo



D	Pevelopment Process	
Review Guidelines Apply to Guidelines Structure practices Review with Group	Pilot Nomenclature Uve in Volunteer Clinics and NRG	Design Process Changes Only minor changes from Pilot
	Slide court (U. Mich.)	esy of Chuck Mayo
ASTRO 2016 ENHAT	NCING ALUE	



#### Task Group findings are in parent committee review process

- Guidelines
  - Target Structures
    - Standardized rule based approach (10) Addresses primary issues and expandable
  - Non-Target Structures
    - Rule based approach (15) with a few concessions
      Specific listing of 736 defined structures
- DVH Nomenclature

#### Slide courtesy of Chuck Mayo

	(U. Mich.)
ASTRO 2016 ENHANCING	ALUE ING OUTCOMES

SILCORE NUME	Patred?	Structure Name	Paired
ANAL CANAL		MAIN BRONC	
BLADDER		OFTIC_NRV	1/.8
BRAC, PLX	_L/_R	ORAL_CAVITY	
BRAIN		OVARY	L/.R
BRAINSTEM		PAROTID	1/ 8
BREAST	L/R	PENILE BULB	
BRONC_TREE		PERINEUM	
CARINA		PELARYNX	
CAUDA_ BOUINA		PITUTTARY	
CEREBELLUM	1/ 8	PROSTATE	
CEREBRUM	L/ B	RECTUM	
CHIASM		RETINA	317.8
CN VII	1/ 8	REP	
CN_VIII	L/ R	SACRUM	
COCHLEA	1/.8	SEM VES	
CORNEA	L/ R	SKIN	
DUODENUM		SM BOWEL	
EAR_MID	1/ R	SPINAL_CORD	
EAR_EXT	_L/_R	STOMACH	
ISOPRAGUS		SUBMIND SALV	1/.8
FEMUR	_1/_R	TIMP_LOBE	1/.8
GLOBE	_1/_R	TESTIS	11.8
GLOTTIS		THYROID	
GREAT_VESS		TM_JOINT	11.8
HEART		TONGUE	
KIDNEY	_1/_R	TRACHEA	
LG_ROWEL		URETHRA	
LARYNX		VULVA	
LAC_GL	_L/_R		
LENS	_L/_R		
LIPS			
LIVER			
LUNG			
A CARLENGER IN			

We're just now agreeing on the ontology of structure names!!!

# Ontology issue example: "How do you designate recurrence in IMRT era?"



-	Type C-low	b	ype B	c	100%
CTVI	Cerebroid Hauter	CTVI CTV2	etra	CTV2	Type E and and a construction of the construct
Type C-int failure	Type A Tailure	Type D-int failure 8. Peripheral high dose failure	Type D-low Type D-low	Peripheral elective dose failure	Стиз
	0.000.00010.000000000000000000000000000				E. Extraneous dose fail
raditional/historic iomenclature	• "Full hit"/"In-field"	"Partial ht"/ "Marginal"	· "in-field"	"Partial hit"/ "Marginal"	E. Extraneous dose fail "Full miss"/ "Out of field"
aditional/historic smenclature entroid location of origin	"Yull Nt"/"to-field"     High dose TV	"Partial hit"/ "Marginal"     High dose TV	"In-field"     Intermediate or low dose TV	"Partial hit"/ "Marginal" • Intermediate or low dose TV •	E. Extraneous dose fai "Full miss"/ "Out of field" Outside all TVs
raditional/historic amenclature antroid location of origin ose to 95% failure volume D95%)	"Yull hit"/"in-field"     High dose TV     295% of the dose     prescribed to TV of origin	"Partial htt"/ "Marginal"     High dose TV        < High dose TV	"In-field"     Intermediate or low dose TV     295% of the dose prescribed to TV of origin	"Partial hit"/ "Marginal"     Intermediate or low dose TV     (95% of the dose prescribed to     TV of origin	E. Extraneous dose fai "Full miss"/ "Out of held" Outside all TVs N/A
raditional/historic amendature antroid location of origin lose to 95% failure volume 095%) usable causes of failure	<ul> <li>"Full het"/"m-field"</li> <li>High dose TV</li> <li>24556 of the dose prescribed to TV of origin</li> <li>Biological failure</li> </ul>	"Partial/htt"/ "Marginal"     High does TV     does TV     does TV     does to TV of origin     to TV of origin     Target defineation error     Dosimetric failure     Doregrown recurrence (paped progressive disease or     neglected late diagnosed recurrence)	Ten field     Intermediate or low dove TV     255% of the dove prescribed     to TV of origin     Improper risk assessment	**Partial htt"/ "Marginal" Intermediate or low does TV • Intermediate or low does TV • Official the does preschool to • TV af cogin • Target deliveration error • Deservice tailore • Deservice tailore • Deservice tailore • Reported Like diagnosed recurrence (************************************	E. Extraneous dose fail "Full miss"/ "Dut of field" Outside all TVs N/A Improper risk assessment Aberrant areas of recurrence

The Standards Problem/Opportunity						
HO' (SEE: A/C CHARC	W STANDARDS PROLIFERAT BERS, CHARACTER ENCODINGS, INSTANT ME	TE: Issaging, etc.)				
SITUATION: THERE ARE 14 COMPETING STANDARDS.	H1?! RIDICULOUS! WE NEED TO DEVELOP ONE UNIVERSAL STANDARD THAT COVERS EVERYONE'S USE CASES. YEAH!	(500N:) SITUATION: THERE ARE 15 COMPETING STANDARDS.				

#### OURNAL OF MEDICAL INTERNET RESEARCH oi:<u>10.2196/jmir.5870</u>

Luo et al

Suidelines for Developing and Reporting Machine Learning Predictive Models in Biomedical Research: A Multidisciplinary







#### JOURNAL OF MEDICAL INTERNET RESEARCH

Guidelines for Developing and Reporting Machine Learning Predictive Models in Biomedical Research: A Multidisciplinary View

7	Prepare data for model building	Identify relevant data sources and quote the ethics approval number for data access.
		State the inclusion and exclusion criteria for data.
		Describe the time span of data and the sample or cohort size.
		Define the observational units on which the response variable and predictor variables are defined.
		Define the predictor variables. Extra caution is needed to prevent information leakage from the re
		sponse variable to predictor variables. <sup>c</sup>
		Describe the data preprocessing performed, including data cleaning and transformation. Remove outliers with impossible or extreme responses; state any criteria used for outlier removal.
		State how missing values were handled.
		Describe the basic statistics of the dataset, particularly of the response variable. These include the ratio of positive to negative classes for a classification problem and the distribution of the respons variable for regression problem.
		Define the model validation strategies. Internal validation is the minimum requirement; external validation should also be performed whenever possible.
		Specify the internal validation strategy. Common methods include random split, time-based split, and patient-based split.
		Define the validation metrics. For regression problems, the normalized root-mean-square error should be used. For classification problems, the metrics should include sensitivity, specificity, pos
		itive predictive value, negative predictive value, area under the ROC <sup>d</sup> curve, and calibration plot [19]. <sup>6</sup>
		For retrospective studies, split the data into a derivation set and a validation set. For prospective studies, define the starting time for validation data collection.

JOURNAL OF MEDICAL INTERNET RESEARCH

Luo et al

Guidelines for Developing and Reporting Machine Learning Predictive Models in Biomedical Research: A Multidisciplinary View doi:10.2196/jmir.5870

Identify and remove redundant independent variables. Identify the independent variables tatum ysuffur from the perfect separation problem. <sup>1</sup> Report the number of independent variables, the number of positive examples, and the number of engative examples. Assess whether sufficient damker of observations in the housive and negative causes. Determine as of candidate modeling techniques (eg. logistic represented, reductions) there also did so a sufficient number of observations in the housive and negative causes. Determine as of candidate modeling techniques (eg. logistic represented, radiation of news) or deg learning). If only one type of model, we used, highly the decision for using that model. <sup>8</sup> Define the performance metrics to select the best model. Specify the model selection, intrapic, Common methods include K-fold validation, proper statification by the response variable is needed. <sup>8</sup> For model selection, include discussion on (1) balance heresem model accurse; and model aligned interpret observations in the observation on the observation in the observation of	8	Build the predictive model	identity independent variables that predominantly take a single value (eg, being zero 99% of the time).
Identify the independent variables that may suffice from the perfect separation problem. <sup>4</sup> Report the number of independent variables, the number of positive examples, and the number of Ansers vehreals efficient data are served indefer for a post of for hered. In particular, the classification data we should be a sufficient tanta are served or of abservations in both positive and negative classes. Determine as set of candidate modeling techniques (e.g., logistic regression, nucleon forest, or deer learning). If only one type of model was used, justify the decision for sing that model. <sup>8</sup> Determine the performance metrics to select the best model. Specify the model selection strategy. Common methods include K-fold validation, or poper stratification by the response variable is needed. <sup>8</sup> For model selection, include discussion on (1) blance hereseen model accuracy and model inger For model selection, include discussion on (1) blance hereseen model accuracy and model inger For model selection, include discussion on (1) blance hereseen model accuracy and model inger For model selection, include discussion on (1) blance hereseen model accuracy and model angel ingo in integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy and model angel is on integration of the selection include discussion on (1) blance hereseen model accuracy a			Identify and remove redundant independent variables.
Report the number of independent variables, the number of positive examples, and the number of negative examples. Assess whether sufficient number of observations in those positive and negative classes. Determine as of candidate modeling techniques (eg., logistic represented classes. Determine): If only one type of models was used, naity the detainion framing that model. <sup>48</sup> Define the performance matrix to select the best model. Specify the model selection array classes, Common methodox include K-fold validation proper statistication on a grid of candidate parameter values. For K-fold validation, proper statistication by the response variable is needed. <sup>38</sup> For model selection, include discussion on (1) bhance hereseen model accuracy and model simplify by or interpretability, and (2) the functions values model accuracy and model simplify by or interpretability.			Identify the independent variables that may suffer from the perfect separation problem. <sup>f</sup>
A sease whether sufficient tanks are available for a good fir of the model. In particular, for classification there should be a sufficient tankness of observations in those pointive and negative classes. Determine a set of candidate modeling techniques (eg., logistic represented eases.) Determine, 14 roly one type of model was used, justify the decision for using that model. <sup>4</sup> Define the performance matrix to select the best model. Support the selection of			Report the number of independent variables, the number of positive examples, and the number of negative examples.
Determine a set of candidate modeling techniques (eg. logitist regression, random forest, or dee learning). If only one type of model was used, havity the detainion for using that model. <sup>8</sup> Define the performance metrics to select the best model. Specify the model selection arranges, Common methodis leadus K-fold validation or boostraps to estimate the lost function on a grid of candidate parameter values. For K-fold validation, proper statistication by the response variable is needed. <sup>8</sup> For model selection, include discussion on (1) balance hereseen model accuracy and model aim ing variables in the selection are selected as a selection of the			Assess whether sufficient data are available for a good fit of the model. In particular, for classification, there should be a sufficient number of observations in both positive and negative classes.
Iteming). If only one type of model was used, justify the decision for using that model. <sup>8</sup> Define the performance metrics to select the best model. Define the performance metrics to select the best model. Specify the model selections strategy. Common methods include K-fold validation, proper estimate the loss function on a grid of candidate parameter values. For K-fold validation, proper stratification by the response variable is needed. <sup>8</sup> For model selection, include discussion on (1) blance between model accuracy and model simpli ity or interpretability, and (2) the function.			Determine a set of candidate modeling techniques (eg, logistic regression, random forest, or deep
Define the performance matrix to select the best model. Specify the model selection strategy. Common methods include K-fold validation or boostraps estimate the lost function on a grid of candidate parameter values. For K-fold validation, proper stratification by the response variable is needed. <sup>3</sup> For model selection, include discussion on (1) balance heresem model accuracy and model simplify by or interpretability, and (2) the diministry with the modeling textingues of the out user. <sup>1</sup>			learning). If only one type of model was used, justify the decision for using that model.8
Specify the model selections transpy. Common methods include K-field validation or bootstrap to estimate the loss function on a grid of candidate parameter values. For K-fold validation, proper stratification by the response variable is needed. <sup>8</sup> For model selection, include discussion on (1) blance between model accuracy and model simpli ity or interpretability, and (2) the minimizity with the modeling techniques of the outset.			Define the performance metrics to select the best model.
stratification by the response variable is needed. <sup>16</sup> For model selection, include discussion on (1) balance between model accuracy and model simpli ity or interpretability, and (2) the familiarity with the modeling techniques of the end user. <sup>1</sup>			Specify the model selection strategy. Common methods include K-fold validation or bootstrap to estimate the lost function on a grid of candidate parameter values. For K-fold validation, proper
For model selection, include discussion on (1) balance between model accuracy and model simpli ity or interpretability, and (2) the familiarity with the modeling techniques of the end user. <sup>1</sup>			stratification by the response variable is needed. <sup>h</sup>
ity or interpretability, and (2) the familiarity with the modeling techniques of the end user. <sup>1</sup>			For model selection, include discussion on (1) balance between model accuracy and model simplic-
			ity or interpretability, and (2) the familiarity with the modeling techniques of the end user. <sup>1</sup>

#### JOURNAL OF MEDICAL INTERNET RESEARCH

Luo et al

Guidelines for Developing and Reporting Machine Learning Predictive Models in Biomedical Research: A Multidisciplinary View

 Table 5. Items to include when reporting predictive models in biomedical research: discussion section.

 Items to include when reporting predictive models in biomedical research: discussion section.

 Items
 Topic

 Or Clinical implications
 Report the clinical implications derived from the obtained predictive performance. For example, report the doll-low range the model area most that could be saved with better prediction. How many patients could benefit from a care model low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model predictive performance. For example, report the doll-low range the model model is not predictive performance. For example, report the doll-low range the model is no fifted that an out that could be refit from a care model is no fifted that model in modeling - Generalizability of the data.

 12
 Unexpected results during the expert
 Report method signed Coefficients, indicating collinearity or complex interaction between method

# Baby steps...

- Why is clinical implementation so hard?
  - Nature of Big Data
    - Sparseness of Big Data
    - Quality assurance
  - Nature of machine learning
    - Algorithms only "know what you tell them, or ask what you ask them to find"
    - Interpretability of complex datasets harder for humans than for algorithms sometimes





# 32 Flavors of machine learning algorithms...

Algorithm	Advantages	Limitations
Decision Tree	<ul> <li>Easy to understand</li> </ul>	<ul> <li>Classes must be mutually exclusive</li> </ul>
	<ul> <li>Fast</li> </ul>	<ul> <li>Results depend on the order of attribute selection</li> </ul>
Carl and a strategy of the		<ul> <li>Risk of overly complex decision trees</li> </ul>
Naive Bayestan	<ul> <li>Easy to understand</li> </ul>	<ul> <li>Variables must be statistically independent</li> </ul>
	<ul> <li>Fast</li> </ul>	<ul> <li>Numeric attributes must follow a normal distribution</li> </ul>
	<ul> <li>No effect of order on training</li> </ul>	<ul> <li>Classes must be mutually exclusive</li> </ul>
		<ul> <li>Less accurate</li> </ul>
k-nearest Neighbors	<ul> <li>Fast and simple</li> </ul>	· Variables with similar attributes will be sorted in the same
	<ul> <li>Tolerant of noise and missing values in data</li> </ul>	<ul> <li>All attributes are equally relevant</li> </ul>
	<ul> <li>Can be used for non-linear classification</li> </ul>	· Requires considerable computer power as number of varial
	<ul> <li>Can be used for both regression and classification</li> </ul>	
Support Vector Machine	<ul> <li>Robust model</li> </ul>	<ul> <li>Slow training</li> </ul>
	<ul> <li>Limits the risk of error</li> </ul>	<ul> <li>Risk of overfitting</li> </ul>
	<ul> <li>Can be used to model non-linear relations</li> </ul>	<ul> <li>Output model is difficult to understand</li> </ul>
Artificial Neural Network	<ul> <li>Tolerant of noise and mission values in data</li> </ul>	<ul> <li>Output model is difficult to understand /s black how a).</li> </ul>
and Deen Learning	<ul> <li>Can be used for classification or memorian</li> </ul>	<ul> <li>Bick of coardition</li> </ul>
and beep tearing	<ul> <li>Can be used for characteristic regression</li> </ul>	- Rak of overlating
	<ul> <li>Can be easily updated with new data</li> </ul>	<ul> <li>Requires a lot of computer power</li> </ul>
		• entropy of the second sec

# Remember, a HUMAN PHYSICIAN has to understand the ML output...



FIGURE 2.7. A representation of the tradeoff between flexibility and interpretability, using different statistical learning methods. In general, as the flexibility of a method increases, its interpretability decreases.

# Leo Breiman, 1928 - 2005

1960 -1967 UCLA (mathematics)

1969 -1982 Consultant

1982 - 1993 Berkeley (statistics)

1984 "Classification & Regression Trees" (with Friedman, Olshen, Stone)

1996 "Bagging" 2001 "Random Forests"



#### RECURSIVE PARTITIONING ANALYSIS (RPA) OF PROGNOSTIC FACTORS IN THREE RADIATION THERAPY ONCOLOGY GROUP (RTOG) BRAIN METASTASES TRIALS

LAURIE GASPAR, M.D., \* CHARLES SCOTT, M.S.,<sup>†</sup> MARVIN ROTMAN, M.D.,<sup>†</sup> Sucha Asbell, M.D.,<sup>†</sup> Theodore Phillips, M.D.,<sup>†</sup> Todd Wasserman, M.D., \* W. Gillies McKenna, M.D., Ph.D. \*\* and Roger Byhardt, M.D.<sup>††</sup>



Fig. 1. Protocol schemas.

Int. J. Radiation Oncology Biol. Phys., Vol. 37, No. 4, pp. 745-751, 1997





# Result: Rad Onc loves RPA!!

- RPA is a "brand name"
  - Classification and regression trees
  - Decision trees
  - Tree-based analysis
- We got into the game early...
  - But we have a long ways to go in machine learning

# Example: ML for toxicity/dose response



logy 118 (2016) 304-314





# RPA-based method

Recursiv	e partition	ning analysis					
Muscle OAR	V-level	Percent- threshold (%)	ROC AUC cohort (test)	ROC AUC holdback (verification)	LogWorth	p-Value	SS
ADM	60	79	0.68	0.60	5.95	<.0001	
BM	35	65.8	0.65	0.57	1.09	0.0815	n.s.
CPM	45	0.35	0.64	0.51	1.00	0.0998	n.s.
GGM	35	98.9	0.70	0.55	2.74	0.0018	
IPC	70	98.2	0.60	0.51	1.08	0.0831	n.s.
ITM	47	99.9	0.67	0.44	2.83	0.0015	
LPM	66	13.1	0.53	0.35	1.07	0.0860	n.s.
LRX	63	1	0.61	0.47	0.89	0.1274	n.s.
MHM	69	17.5	0.74	0.64	6.77	<.0001	
MM	66	4.4	0.61	0.53	0.88	0.1314	n.s.
MPC	49	99.9	0.63	0.54	0.17	0.6825	n.s.
MPM	70	1	0.59	0.45	3.31	0.0005	
PDM	69	13.5	0.60	0.48	0.15	0.7070	n.s.
PGM	65	68.9	0.62	0.49	0.24	0.5732	n.s.
SPC	70	6.35	0.68	0.47	5.09	<.0001	

Statistically significant at P < 0.05. Statistically significant after Bonferroni correction.

			RI	PA-	bas	sed n	netł	nod			
		Model	Whole		ROC	POC			Effect		
Model	Model effects	good- ness of fit*	model LogWorth †	Whole model p-value	AUC Cohort (test)	AUC Holdback (verification)	Model BIC [ΔBIC]	Evidence Grade §	likelihood LogWorth †	Effect likelihood p-value	Odds Ratio (95% CI)
0	MHM V69			- 0001		0.032	124.25	BIC	5.19	<.0001	1.09 (1.03-1.16) ‡
Continuous	Age	0.9999	0.94		0.835	0.812	[0]	(reference)	2.56	0.0027	1.03 (1.02-1.05) ‡
Disas	MHM V69 >79.5%	0.9997	8.76	< 0001	0.825	0.864	130.53	Steven	7.62	<.0001	4.8x10 <sup>7</sup> (12.63-∞) ◊
	Age >62					0.004	[6.28] §	Same	1.37	0.043	2.57 (1.03-6.53)

MDACC RPA-based method showed Myelohyoid muscle V69 and age were best predictors...

#### E-idence Gradia ( stational stationa



Development of a multivariable normal tissue complication probability Cross? (NTCP) model for tube feeding dependence after curative radiotherapy/ chemo-radiotherapy in head and neck cancer

Kim Wopken<sup>a,\*</sup>, Hendrik P. Bijl<sup>a</sup>, Arjen van der Schaaf<sup>4</sup>, Hans Paul van der Laan<sup>4</sup>, Olga Chouvalova<sup>4</sup>, Reel J.H.M. Steenbakkers<sup>4</sup>, Patricia Doornaert<sup>9</sup>, Ben J. Slotman<sup>5</sup>, Sjoukje F. Oosting<sup>5</sup>, Miranda E.M.C. Christiane<sup>4</sup>, Bernard F.A.M. van der Laan<sup>4</sup>, Jan L.N. Roodenburg<sup>4</sup>, C. René Leemans<sup>4</sup>, Irma M. Verdonck-de Leeuw<sup>4</sup>, Johannes A. Langendijk<sup>4</sup>

/ariable	OR	OR 95% CI	p-Value
-classification			
Tis-T2	1.00		
T3-T4	1.53	(1.17 - 2.06)	< 0.001
Baseline weight loss			
No weight loss	1.00		
Moderate weight loss (1-10%)	2.58	(2.01 - 3.19)	<0.001
Severe weight loss (>10%)	5.08	(3.32 - 7.30)	<0.001
Treatment modality			
Conventional fractionation	1.00		
Radiotherapy + cetuximab	1.74	(1.50 - 2.01)	< 0.001
Accelerated fractionation	3.33	(2.40 - 4.53)	< 0.001
Chemoradiation	6.73	(4.00 - 10.98)	<0.001
Dosimetric variables			
PCM superior mean dose (Gy)	1.07	(1.04 - 1.09)	<0.001
PCM inferior mean dose (Gy)	1.03	(1.01 - 1.05)	0.006
Contralateral parotid mean dose (Gy)	1.01	(1.00 - 1.02)	0.14
Cricopharyngeal muscle mean dose (Gy)	1.02	(1.01 - 1.03)	0.004

**UMCG** Lasso method showed pharyngeal constrictor muscle and cricopharyngeu s were best predictors... Radiotherapy and Oncology 113 (2014) 95-101

# Who was right?

- · Answer: Everyone!
  - Solid methods/different datasets
  - Dose variation was due to tumor location/site, and dose to OARs, not "magic dose threshold"
  - MDACC all oropharynx, block over cricopharyngeus m., posterior constrictors m. with low variability as they were partially covered by RP node CTVs.
  - UMCG mostly larynx, little dose to myelohyoid m.

5



Fig. 1 t of the o elation coefficients among predictor variables nt colors from green to red. Red indicates high two variables. Here variables are named by ir description is listed in Table.



Fig. 2. A box plot of the fit performance in training sets an prediction performance in test sets for xerostomia data of all variable model is shown. The median of the likelihood is acted by the center line, and the fits and third quariles are the edges of the box area. The extreme values (within 15 times thinterquarile narge from the upper of lower quarile) are the end of the lines extending from the interquarili range. Points a distance from the median grate than 1.5 times the interquarili range are plotted individually (circles).

Xu et al. Int J Radiation Oncol Biol Phys, Vol. 82, No. 4, pp. e677-e684, 2012

# Problem/opportunity: We don't have accessible aggregate data...





# Peeling the onion...



FIGURE 7.1 Modern radiation oncology practice comprises five major components: peopleware, processes, information, software, and hardware. Successful development, implementation, and use of IS and IT applications requires proper consideration of the hierarchy as well as synergy among each of the components.







Figure 1. Chart illustrates
the radiation therapy work
flow for a patient with prostate
cancer. The work flow consists
of two stages, treatment plan-
ning (steps 1-18) and treat-
ment delivery (steps 19-30),
and is broken down into steps
that can be used to evaluate
the efficiency of the system
once it is implemented. The
work flow starts when the ra-
diation oncologist decides at
consultation that the patient is
to receive radiation treatment
(step 2). Gray boxes indicate
steps that are performed by
the radiation oncologist. Of
these steps, steps 9 and 13 are
incorporated into the treat-
ment planning system (TPS)
in most cases, but steps 15,
24, 27, and 30 either require
hard-copy records or radia-
tion therapy information that
may not be readily accessible.
DRR = digitally reconstructed
radiograph, DVH = dose-
volume histogram, LINAC =
linear accelerator, QA = qual-
ity accurance























#### The Cancer Journal • Volume 17, Number 4, July/August 2011

- TABLE 1. Producers of Information in Clinical Trials
- Producers of Information in Clinical Trials Are as Follows:
- A. Pre-enrollment The authors of the protocol
- The IRBs that approve it for local participation
   The lead institution, sponsor, or CRO that creates the CRFs for EDC
- CRFs for EDC The individuals (coordinators/clinicians, individual investigators, and budgeting/billing administrators) that interpret data for that protocol into local systems, including scheduling, budgeting, reviews, and departmental and institutional reporting. Clinical trial budgeting and contract management staff
- B. Intraenrollment The subject/patient
- In a subject/patient
   The clinician, both in CRFs and in source documentation
   (the medical record—on paper or electronic)
   The laboratory (hematology, chemistry genomics,
   microbiology, pathology) and radiology departments
   and systems
   Research memister (medicate and medicate)
- Research monitors (queries and reviews)
- Internal and external regulatory reviewers
   Internal and external auditors



FIGURE 5. Validation for format, fields, and values against standards: a simple configuration for standards designers.









# DICOM-RT and Its Utilization in Radiation Therapy $^{\rm i}$

#### Maria Y.Y. Law, PhD \* Brent Liu, PhD



Information Entity of DICOM Image Object	Module* (MR Imaging)	Information Entity of DICOM-RT Object	Module* (RT Structure Set)
Patient	Patient	Patient	Patient
Study	General Study	Study	General Study
Series	General Series	Series	RT Series
Frame of Reference	Frame of Reference		
Equipment	General Equipment	Equipment	General Equipment
Image	General Image, Imaging Plane, Image Pixel, MR Image, SOP Common (13)	Structure Set	Structure Set, ROI Contour RT ROI Observations, SOP Common (13)

tient's Sex, and so on. The modules for the image Information Entity depend on the modality concerned. In this table, the modality of interest is MR imaging; thus; in addition to the General Image module, Information Entity has a specific MR Image module and other related modules. Structure Set is one of the DICOM-RT objects and includes Structures Set, ROI Contour, and RT ROI Observations as its specific modules. IOD = information object definition, ROI = region of interest, SOP = service-object pair.







FIGURE 7.5 Hardware components and connectivity infrastructure of the modern radiation oncology practice. Highly available data storage and efficient information flow are highly dependent on having and maintaining reliable, safe, and stable infrastructure and connectivity across all the systems inside and outde the radiation oncology department.















Figure 3 Integrating imaging and clicical information to construct a disease model. (A) Features are computed from normalized imaging taudies, in this case, a lang cancer is identified on a CT imaging study top imagel, and a corresponding position emission tomography study (bottom image) is performed as assess the level of nonweractivity. (B) An Integrated data model creates the basis for combining image data with additional information from the patient's medical record, thereby providing context to the imaging study integrated start model is in individual's starts. (C) A patient's smooting history, results of thoracity turno Elopsies, and course (and a patient is model and in an individual's starts. (C) A patient's moding history, results of thoracity turno Elopsies, and course (and a patient is model and in an individual's starts. (C) A patient's moding history even individual, of there gained model and be used as propositios to for a given individual, of there gained model has been writed for a given population, them he model can be used as be embedded writhin the disease model.

1056

Bui AAT, et al. J Am Med Inform Assoc 2013;20:1053-1058. doi:10.1136/amiajnl-2012-001340













Key element. category	Demand ranking	ETL difficulty	Typical source systems	Access	Mahiple source systems	Use or used free text entry	Missing data	Data accuracy	Lack of standar dization	PHI constraints Ernit access	Legacy formats or systems	Require process changes	Extensive transformation	Other
Demographics	1	L	EHR	×										E
Health status factors	2	L	EHR	×										B
Pathology ()	3	M to H	EHR	×		×	×		×		×	8		E, X
Surgery O	2	M to H	EHR	×		×	×		×		×	8		E, X
Chemotherapy	2	М	EHR, ODB	×										Е
Encounter details  Office, emergency room, hospitalization	3	L	EHR	8									×	R
Diagnosis 🔍 , , +	1	м	EHR, ROIS	×	×			×			×	8		R, E
Staging 0,+,+	1	н	EHR, ROIS	×	×	×		×			×	8		Е
Prescription +, •	1	н	ROIS,						8			×		E, X, R
Astronted	1	м	ROIS										×	
nim details														
DVH . n .	1	м	TPS				×		× .		×	8	×	ATPS
Sarvival	1	м	EHR, XLS, ODB	×						8				UD, E
Recurrence +.+	1	н	EHR	×		×	×			×	×	8		E X
Toxicity	1	н	EHR, ROIS	×		×	×			×	×	8		E, X
Patient-reported outcomes +	2	н	EHR, P	×			×			×	×	8		Е, Х
Laboratory values	2	м	EHR	8				×					×	Е
Medications	2	M	EHR	8				×					×	Е
Height, weight,	2	м	EHR	8				×					×	Е
Treatment imaging: Timeline details	3	н	ROIS										×	R
Diagnostic imaging details ()	3	М	ODB	8	×						×		×	
Rediomics O.+	3	L	XLS	×									8	
Generation (1)	3	L	XLS	×									8	
Charas .	3	1.	ROIS											
Research	4	н	XLS					×	8			×	×	Е
Constant of Consta			ODB											



# What about outcome/clinical/PRO integration?

#### Can electronic web-based technology improve quality of life data collection? Analysis of Radiation Therapy Oncology Group 0828

Benjamin Movsas MD<sup>a,\*</sup>, Daniel Hunt PhD<sup>b</sup>, Deborah Watkins-Bruner PhD, RN<sup>c</sup>, W. Robert Lee MD, Med<sup>d</sup>, Heather Tharpe RN OCN<sup>c</sup>, Desiree Goldstein RN MSN<sup>f</sup>, Joan Moore RN OCN<sup>g</sup>, Ian S. Dayes MD<sup>b</sup>, Sara Parise RN OCN<sup>f</sup>, Howard Sandler MD<sup>j</sup>

> The EPIC QOL compliance rates at baseline, 6 months (the primary endpoint), and 12 months are shown in Table 2. At baseline, the QOL completion rate was 98%. Compared with the 52% 6-month EPIC completion rate using paper forms (in RTOG-0415), the EPIC QOL web-based completion rate at 6 months was 99% (2-sided P value < 001). At 1 year, the EPIC QOL web-based compliance rate was 82%. Reasons for noncompliance at 6 months with EPIC QOL via the web-based strategy were patient relias(12%), patient could not be contacted (2%), or other reason (6%).

#### Computerized patient-reported symptom assessment in radiotherapy: a pilot randomized, controlled trial

Erik K. Fromme<sup>1</sup> - Emma B. Holliday<sup>2</sup> - Lillian M. Nail<sup>3</sup> - Karen S. Lyons<sup>3</sup> -Michelle R. Hribar<sup>4</sup> - Charles R. Thomas Jr<sup>5</sup> Support Care Concer (2010) 24:1877–1906 DOI 10.1007/a00529-015-288-3



NALE OF TAXABLE PARTY O		Excluded (n = 53) Nor maximg inclusion orbits (c) = 80 References for participue (n = 20) Other measure (n = 5)
	Bankerised (n	-112)
Almados	Allocated to Control (2 = 24). Reserved idecond intervention (a = 54). Did on receive discand	Allocated to Intervention (\$ = 2.8) Bounted allocated intervention (\$ = 5.9) Did not reaction disconted
Fullow up	Less to follow up (ti - 0) Discontinued movember (k - 2) Meaning respective (to -1)- Protect conserve withdrawe (k-2)	Lore to follow up (n ~ 6) Discolutional intervision (n ~ 1) Reductor sugged due to disease progenities (n ~ 1)
when	Analyzed (n = 30). Enduded from analysis (n = 0). Printing did not have at loast 4 complete analysis.	Analyzed (z = 57) Escluted from analysis (a = 1) Patiente dat net have at loar 4 completed successors. Sor publics



Fig. 3 Sample report showing what was printed and given to clinicians in the intervention group. Patients received a version with larger font occupying two pages front and back. The *np* shows the survey dates and columns from the lift to the right show how scores change over time. The first show rows show scores for Kamofsky Performance Shun (KPS), Quality of Lift (QOA), and satisfaction with pair control shun (KPS).

The next stree rows show scores for pain in different automic locations. The next state rows show Mamcrial Symptom Assessment Scale (MSAS) scores. The following tor rows show skin truckity symptoms, other sizespecific physical symptoms, and emotional and function scales. The bottom rows give a key to interpret the patient-reported outcomes (PROs) for symptoms, F/S, and OL/pain corted suitsfaction



### Impact of Statistical Learning Methods on the Predictive Power of Multivariate Normal Tissue Complication Probability Models

Cheng-Jian Xu, Ph.D., Arjen van der Schaaf, Ph.D., Cornelis Schilstra, Ph.D., Johannes A. Lannandiik, M.D., Ph.D., and Aart A. van't Vald. Ph.D.

/ariable index	Description	Range or limits	Median or frequency	Correlation
1	Chemotherapy*	0, 1	142, 43	0.1905
2	Gender	0, 1	118, 67	0.0316
3	Age	4092	62	0.0418
4	Medical center <sup>1</sup>	0, 1	140, 45	0.1983
5	Volume of soft palate	1.5-12	4.56	-0.1184
6	Mean dose to soft palate (Gy)	0-73	25.96	0.4123
7	Volume of contralateral parotid gland	9-58	25.10	-0.1360
8	Mean dose to contralateral parotid glands (Gy)	0-69	26.37	0.4935
9	Volume of ipsilateral parotid gland	9-56	26.31	-0.1356
10	Mean dose to ipsilateral parotid gland(Gy)	0-71	35.16	0.4634
11	Volume of contralateral sublingual gland	0.06-2.5	0.39	-0.0571
12	Mean dose to contralateral sublingual gland (Gy)	0-73	14.53	0.3070
13	Volume of ipsilateral sublingual gland	0.06-2.4	0.35	-0.1355
14	Mean dose to ipsilateral sublingual gland (Gy)	0-73	16.31	0.2574
15	Volume of contralateral submandibular gland	3.5-17	8.65	-0.1228
16	Mean dose to contralateral submandibular gland (Gy)	0.05-76	52.70	0.4532
17	Volume of ipsilateral submandibular gland	0.6-20	8.71	-0.0886
18	Mean dose to ipsilateral submandibular gland (Gy)	0.0676	54.71	0.4362
19	Volume of the lower lip	0.5-5.5	1.98	-0.0013
20	Mean dose to the lower lip(Gy)	0-71	0.79	0.1662
21	Baseline xerostomia score(s)	1, 2, 3, 4	106, 53, 19, 7	0.3578

#### BIG DATA

Lessons From Large-Scale Collection of Patient-Reported Outcomes: Implications for Big Data Aggregation and Analytics Jeff A. Sloan, PhD,<sup>+</sup> Michele Halyard, MD,<sup>+</sup> Issam El Naga, PhD,<sup>+</sup> and Charles Mayo, PhD<sup>+</sup>









Patient Reported Outcomes (PROs) in Clinical Trials: Is 'In-Trial' Guidance Lacking? A Systematic Review Derek G. Kyte<sup>1</sup>\*, Heather Draper<sup>2</sup>, Jonathan Ives<sup>2</sup>, Clive Liles<sup>3</sup>, Adrian Gheorghe<sup>1</sup>, Melanie Calvert<sup>1,4</sup> PRE-TRIAL GUIDELINE IN-TRIAL GUIDELINE POST-TRIAL GUIDELINE FUTURE RESEARCH OTHER Т OM EVALUATIONSELECTION STUDY DESIGN & PROCEDURE PROVISION & UNDERSTANDING DATA ANALYSIS/REPORTING/ PRESENTATION LITY CONTROL/COMPLIANCE & CORRECT USE OF OM DATA INTERPRETATION/ LABELLING & PROMOTIONAL CLAIMS VELOPMENT/VALIE MODIFICATION HELP/PROXY ASSES MENTS OTHER REPORTING OF DATA COLLECTION/ SCORING DEALING WITH CONCERNING PRO DATA -JLD PRO DATA INFORM MANAGEMENT PRO DATA ACCESS 

OTHER

OPEN @ ACCESS Freely a

W PLOS

Information type	Data examples	IT syst
Baseline clinical data	Demographics (including co-morbidity and family history), TNM-stage, date of diagnosis, histopathology	HIS, TDS
Diagnostic imaging data	Diagnostic CT, MR and PET imaging	PAC
Radiotherapy treatment planning data	Delineation/structure sets, planning-CT, dose matrix, beam set-up, prescribed dose and fractions	PAC: RIS
Radiotherapy treatment delivery data	Cone beam CTs, orthogonal EPID imaging, delivered fractions	PACS
Non-radiotherapy treatment data	Surgery, chemotherapy	HIS, TDS
Outcome data	Survival, local control, distant failure, toxicity (including patient reported outcomes), quality of life	EDC, TDS
Follow-up imaging data	Follow-up CT, MR and PET imaging	PAC
Biological data	Sample storage, shipping, tracing and lab results	LIMS
Additional study conduct data	Study design, protocol, eligibility criteria	EDC, CTM



#### Data collection

Benefits of a clinical data warehouse with data mining tools to collect data for a radiotherapy trial

Erik Roelofs <sup>a,1</sup>, Lucas Persoon <sup>a,1</sup>, Sebastiaan Nijsten <sup>a</sup>, Wolfgang Wiessler <sup>b</sup>, André Dekker <sup>a,1</sup>, Philippe Lambin <sup>a,1</sup>





© 2012 Ted Goff

you to condense them down to one meaningful warehouse."

#### Data mining for radiotherapy trials 176

Parameter	NSCLC	Rectum	Source	Source		
			Manual	Automatic		
Gender	~	1	Chart	EMR	Looked up	
WHO score	V	V	Chart	EMR		
INM staging	1	4	Chart	EMR		
Chemo therapy	V	4	Chart	EMR		
fumber of positive lymph nodes	V	2	Chart	EMR		
fumour PA	ý.	ý	Chart	EMR		
CR		ý	Chart	EMR		
iurvival	~	2	Chart	EMR		
lotal delivered dose	4	ý	R&V	R&V		
Overall treatment time	V	ý.	R&V	RBV		
GTV volume	1	1	XiO	PACS	Recalculate	
<i>v</i> <sub>1</sub>	Lungs"		XiO	PACS		
(m	Lungs		XIO	PACS		
leo		Bladder	XiO	PACS		
dLD	×		XIO	PACS		
SUV Max	1 * A	Tumour	TrueD	PACS		
SUV Mean		Tumour	TrueD	PACS		

 $^6$  V<sub>3</sub> and V<sub>20</sub> data for the lungs were calculated with both lungs minus the PTV.  $^6$  MLD data for the lungs were calculated with both lungs minus the GTV.





# Radiomics/radiogenomics magnifies this challenge/opportunity

- Radiomics is "extraction and analysis of large amounts of advanced quantitative imaging features with high throughput from medical images obtained with computed tomography, positron emission tomography or magnetic resonance imaging." (Kumar, 2012)
   Radiogenomics

Imaging

- In imaging community refers to linkage of imaging features w genomics · In radiation oncology community refers to genomic correlates of radiation
- response

Fig. 4. The Radiomics workflow. On the medical images, segmentation is performed to define the tamour region. From this region the features are extracted, e.g. features based on tamour intensity, texture and shape. Finally, these features are used for analysis, e.g. the features are assessed for their proposite power, or linked with size, or gene expression.

Feature extraction

Analysis

Exciting time!!

Segmentation









Lung Texture in Serial Thoracic Computed Tomography Scans: Correlation of Radiomics-based Features With Radiation Therapy Dose and Radiation Pneumonitis Development Alexandra Cunliffe, PhD, \* Samuel G. Armato III, PhD, \* Richard Castillo, PhD, <sup>1</sup>Ngo PhD, <sup>1</sup>Ngo PhD, <sup>1</sup>Samuel G. Armato III, PhD, <sup>1</sup> and Hania A. Al-Hallaq, PhD





Fig. 2. Plots of  $\Delta FV$  versus mean physical dose with 95% confidence intervals for 2 features. Average values were calculated on a per-patient basis in 5-Gy dose bins.

Table 4 Comparison of single-feature performance to distinguish patients with RP compared with a classifier composed of each 2- or

No. of features/dose-dependent	Mean AUC (range)							
measurement	Low dose	Medium dose	High dose	Fitted slope				
1 Feature	0.64 (0.49-0.72)	0.68 (0.49-0.77)	0.71 (0.52-0.78)	0.71 (0.50-0.78)				
2 Features	0.66 (0.59-0.74)	0.71 (0.59-0.78)	0.73 (0.59-0.78)	0.74 (0.59-0.84)				
3 Features	0.66 (0.59-0.75)	0.72 (0.59-0.79)	0.72 (0.60-0.79)	0.75 (0.61-0.83)				

#### Conclusions

This study demonstrated that quantitative measurement of dose-dependent texture changes between pre- and post-RT CT scans can differentiate between patients with clinical (grade 23) RY and hose patients without RF. Twelve in-tensity- and texture-based features demonstrated signif-canly increased changes for patients with RF in general, individual features could be used to discriminate between patients with and without RF with moderate performance. When multiple features were combined in a classifier, AUC increased significantly (AUC values from 0.59 to 0.84). This study demonstrates the potential ability of radiomics provide a quantitative, individualized approach to mea-sure patient lung tissue reaction to KT and assess radiation pneumonitis. pneumonitis.













However, when GADD34 was induced and radiation therapy was given, there was additional growth delay due to a decreased hypoxic fraction at the time of irradiation, as previously described [12]. Thus the difference in tumor heterogeneity between gene-induced and non-induced tumors post RT was reflected in the imaging features likely as a result of a phenotypic change. Interestingly, the radiomics image features that were found to be significantly different between both groups shortly after irradiation were also observed at larger tumor volumes. This phe-nomenon was observed independent of the CT image acquisition energy level, although the observed features were different between both energies tested. Remarkably, the feature value for slow-growing tumors (gene-induced) was higher than for faster-growing tumors (no gene-induced group) upon combination with radiotherapy.

#### Conclusion

We have shown in in vivo preclinical models that radiomics is able to quantify the early effects of radiation treatment and genetic changes in tumors with similar volumes, and identify differences that are not visible to the human eye.



Figure 2. (A) Con





# Watch out for the trough!!





# Example: Oropharynx Machine Learning/Radiomics Challenge

- ~20,000 US cases
- Oropharynx cancer (OPC) is epidemic due to human papilloma virus (HPV) infection
- More HPV-associated head and neck cancers than cervical cancers now

	ESTI	ESTIMATED DEATHS				
	BOTH SEXES	MALE	FEMALE	BOTH SEXES	MALE	FEMALE
All Sites	1,529,560	789,620	739,940	569,490	299,200	270,290
Oral cavity & pharynx	36,540	25,420	11,120	7,880	5,430	2,450
Tonque	10,990	7,690	3,300	1,990	1,300	690
Mouth	10.840	6,430	4,410	1,830	1,140	690
Pharynx	12,660	9,880	2,780	2,410	1,730	680
Other oral cavity	2,050	1,420	630	1,650	1,260	390

# Human Papillomavirus (HPV)



- DNA virus
- >100 different sub-types
- Infects skin and mucosa
- Asymptomatic
- Benign growths warts
- Oncogenic (cancer causing) types are mostly 16 and 18







Figure 1 The "radiomics workflow" involves a series of iterative steps for reproducible and consistent extraction of imaging data. These steps include image sequisition, namor segmentation, feature extraction, and feature selection. The selected features can then be analyzed for outcome correlation and potential incorporation into potenticive models. Additionally, validational should be done against completely independent large datasets, preferably from other institutions. Tranul Cancer Res 2016;55(4):371-382

Radiomics	in	head	and	neck	cancers:
	C	lassif	icati	nn	

Table 1 Studies on radiomics for segmentation and classification						
Authors (study)	Publication date	Modality	# of patients	Anatomic site, if specified	Analyzed endpoint	
Raja et al. (23)	Sep 2012	СТ	21	Oral cavity	Tumor grade classification	
Buch et al. (17)	Jul 2015	СТ	40	Oropharyrox	HPV status	
Fujita et al. (24)	Jan 2016	СТ	46	Oropharynx (25); larynx (17); hypopharynx (5)	HPV status	
Yu et al. (26)	Mar 2009	FDG-PET/CT	20	-	Normal vs. abnormal tissue classification	
Yu et al. (27)	Oct 2009	FDG-PET/CT	10	Oropharynx and nasopharynx	Normal vs. abnormal tissue classification	
Vallieres et al. (25)	Oct 2013	FDG-PET	67	-	HPV status; loco-regional failure; distant metastasis	
Fruehweld- Pallamar et al. (28)	Nov 2013	MRI	38	Parotid	Benign vs. malignant status; tumor type differentiation	
Yang et al. (29)	Dec 2014	MRI	15	Parotid	Parotid vs. surrounding tissue differentiation	
Brown et al. (15)	May 2015	DW-MRI	26 (training) 18 (validation)	Thyroid	Thyroid nodule classification	
Jansen et al. (21)	Jan 2016	DCE-MRI	19	Oropharynx	Local control; local failure	
Park et al. (30)	Feb 2016	MRI	27	Oropharynx	Tumor type differentiation	
Fruehwald- Pallamar et al. (31)	Feb 2016	MRI	100	-	Benign vs. malignant status	



# Head and neck radiomics: Outcomes

			Wong et al. Radiomics in head and neck cancer: from exploration to application				
Table 2 Studies on radiomics for prognostic and predictive biomarkers							
Authors (study) Put dat	blication te	Modality	# of patients	Anatomic site, if specified	Analyzed endpoint		
Zhang et al. (22) Des	ic 2013	СТ	72	Oral cavity (28); larynx (21); hypopharynx (14); oral cavity (8)	Overall survival		
Aerts et al. (34) Jur	n 2014	СТ	474 (training); 545 (validation)	Lung or head and neck	Median survival		
Parmar et al. (35) Jur	n 2015	ст	878	Lung or head and neck	Survival; tumor stage; HPV status		
Parmar et al. (36) Des	c 2015	ст	101 (training); 95 (validation)	-	Overall survival		
Leijenaar et al. (37) Aug	ig 2015	ст	542	Oropharynx	Median survival		
El Naqa et al. (20) Jur	n 2009	FDG-PET	9	-	Overall survival		
Dang et al. (38) Jan	n 2015	MRI	16	Oropharynx	p53 status		

#### Transl Cancer Res 2016;5(4):371-382

# Our challenges...

- If image features in the local region can predict local phenotypic alteration (per Panth et al.), can we predict local oncologic outcomes (i.e. local control in the priamry tumor)?
- If the major driver of local control probability genomically is known, can we identify genetic correlates of local control risk (i.e. HPV status)?

# Unmet needs, unanswered questions

- · Can we crowdsource innovation?
- · Can we draw on "non-insider" knowledge?
- · Can we have some fun?



# 2016 OPC CPM challenges

Determine from CT data whether a tumor will be controlled by definitive radiation therapy.

Predict from CT data the HPV phenotype of oropharynx tumors; compare to groundtruth results previously obtained by p16 or HPV testing.

# Kaggle/MICCAI challenges



# Rules of the game

- · Binary classification
  - HPV status (1/0)
  - Local control of tumor in the primary site (1/0)
- Can leverage demographic data
- Half dataset with known HPV/local control available for training
  - Other half randomly split into "public" validation set (visible to end users) and "private" validation set
  - Could submit and see performance on "public" set to avoid overfitting
  - Winner declared on "private" dataset





Additionally, 1st, 2nd, and 3rd place teams will receive a commemorative trophy, ideally of limited cash value and likely with hideous design, allowing gentle flaunting of their definitive victory over anguished peers and colleagues who failed to place in the challenge. An example is shown below.





Results								
Local Control Challenge	HPV Challenge							
1 4 teams	9 teams							
2 3 players	1 7 players							
7 8 entries	7 O entries							



Private Leaderboard - Oropharynx Cancer (OPC) Radiomics Challenge :: Human Papilloma Virus (HPV) Status Prediction

	Δrank	Team Name	Score @	Entries	Last Submission UTC (test - Last Submission)
1	n	BIG-S2_Veera_HPV all KabdanYu Saamine Yang Rongite Lu Yengfei youyi YangYu	0.91549	5	Mon, 12 Sep 2016 23:35:54 (-2.2d)
2	† <b>2</b>	JA st. • johnwshumway • AlokanandaGhosh	0.80047	6	Mon, 12 Sep 2016 22:44:40 (-2d)
3	-	Nguyen Khanh	0.74883	8	Sun, 11 Sep 2016 18:17:31 (-27.3d)
4	11	الحوت الأبيض	0.69190	5	Mon, 12 Sep 2016 20:08:10 (-14.7d)
5	12	JunlinYang	0.67254	17	Thu, 08 Sep 2016 01:53:02 (-10.7d)
6	12	tyler	0.66491	5	Thu, 08 Sep 2016 21:19:17 (-10.6d)
7	11	USF-Moffitt	0.65200	13	Mon, 12 Sep 2016 22:08:02 (-8d)
8	11	turingcomplete	0.58275	6	Sat, 10 Sep 2016 10:12:53
9	,a	The_Courtyard II ccc1986 Rachel xjfave	0.52054	5	Tue, 30 Aug 2016 22:40:25 (-6d)

#### Oropharynx Cancer (OPC) Radiomics Challenge :: Local Recurrence Prediction

	Arach	Team Name	Score @	Detries	Last Submitsion UTC (see - Last television)	
1	12	Nguyen Khanh	0.92405	4	Sun, 11 Sep 2016 18:09:25 (-30.4d)	
2	11	الموث الأبيخن	0.91930	6	Moy, 12 Sep 2016 17:07:23 (-14.8d)	
3	.0	turingcomplete	0.90506	14	Sat, 10 Sep 2016 19:15:55 (-4.1d)	
4	11	JA ili • johneshamsey • AlokanendaGhoah	0.86709	4	Mon, 12 Sep (2016 04:52-64 (-0,34)	
5	14	JunlinYang	0.80538	1	Sun, 28 Aug 2016 07,25c41	
6	11	USF-Moffitt	0.71203	14	Mon, 12 Sep 2016 22:17:21 (-3.10)	
7	2	BIG-S2_Veera_LRP # KatalasiYu Jaamine Yang Tangfel chaohaang.atat yoouj Yangfte	0.69620	12	Moru 12 Sep 2016 23:48:50 (-6.46)	
8	12	ECELLWARRIORSIIH	0.68671	4	Weld, 31 Aug 2016 02:20:14 (-0.110	
9	.1	Arjun	0.67405	3	Wed, 31 Aug 2016 02;2(o44 (-9.2h)	
10	ų	UniNA - GatelenePlantadose - Stefano	0.66772	7	Mon, 12 Sep 2016 23:58:07	
- 11	.5	NinoArsov	0.50000	3	Set, 13 Aug 2016 00:09:13 (-0.2h)	
12		Courtyard d Hachel xylave cec1985	0.50000	2	Mor, 12 Sep 2016 21:28:31 (-38)	
13		albert	0.50000	2	Mon, 12 Sep 2016 04:44:02 (-01)	

•			
-			

# Congrats to our winners!!

1st Place and recipient of a proffered manuscript acceptance (post-editorial review), with waiver of fees in the well-renowned international, open access ESTRO-sponsored journal: Clinical and Translational Radiation Oncology (rEO);

Nguyen Khanh This college kid beat everyone else by \*NOT USING RADIOMICS\* Congrats Nguyen!! But The best is yet to come. You're invited to attend the 19th MICCAI international conference, taking place in Athens, Greece during the period Cottober 17:21, 2016 in person (no travel support will be provided) or via Skype to present your winning algorithm and be celebrated among your fellow data scientists

1st Place and recipient of a proffered manuscript acceptance (post-editorial review), with waiver of fees in the well-renowned international, open access ESTRO-sponsored journal; Clinical and Translational Radiation Oncology (rEO);

#### BIG-S2\_Veera\_HPV

Congratsill Burn The best is yet to come. You're invited to attend the 19th MICCAI international conference, taking place in Athens, Greece during the period October 17-21, 2016 in person (no travel support will be provided) or via Skype to present your winning algorithm and be celebrated among your fellow data scientists.

# TRUTH: Good models are hard to find...so showing your work helps :)

Machine Learning methods for Quantitative Radiomic Biomarkers Chitan Permitty, Patick Grosman<sup>44</sup>, John Bunist, Philos Lambia &

Classification method acronym	Classification method name	Feature Selection method acronym	Feature selection method name Relief				
Nnet	Neural network	RELF					
DT	Decision Tree	FSCR	Fisher score				
BST	Boosting	GINI	Gini index				
BY	Bayesian	CHSQ	Chi-square score				
BAG	Bagging	JMI	Joint mutual information				
RF	Random Forset	CIFE	Conditional informax feature extraction				
MARS	Multi adaptive regression splines	DISR	Double input symmetric relevant				
SVM	Support vector machines	ort vector machines MIM Mutual information					
DA	Discriminant analysis	CMIM	Conditional mutual information maximization				
NN	Neirest neighbour ICAP		Interaction capping				
GLM	Generalized linear models	TSCR	T-test score				
PLSR	Partial least squares and prinicipal componenet regression	MRMR	Minimum redundancy maximum relevance				
_	-	MIFS	Mutual information feature selection				
-	-	WLCX	Wilcown				

												RELF
												FSCR
												GINI
												CHS
												JMI
												CIFE
												DISR
												MIN
												CMI
												ICAP
												TSCF
												MR
				**								MIR
												WLC
Nnet	DT	BST	RY BY	BAG	22	MARS	SVM	đ	NN	GLM	PLSR	
				Class	Beat	inn N	Intho	4.				

# Share your models/approaches!

- Github
  - Software repository
- Figshare
  - Data repository
- The Cancer Imaging Archive
  - Image repository
- · Nature Scientific Data
  - Publish data/software descriptor on PubMed

#### Data Descriptor: Matched computed tomography segmentation and demographic data for oropharyngeal cancer radiomics challenges

SCIENTIFIC DATA | 4:170077 | DOI: 10.1038/sdata.2017.77

Data category	Description					
Patient ID	Numbers_given_randomly_to_the_patient_after_anonymizing_the_DICOM_PHI_tag_(0010.0020;_Patient_ID)					
HPV/p16 status	Human Papilloma Virus status, as assessed by HPV DNA in situ hybridization <sup>11</sup> and/or p16 protein expression via immunohistochemistry (IHC), with the results described as: 1 (i.e., Positive) or 0 (i.e., Negative)					
Gender	Patient's sex					
Age at diagnosis	Patient's age in years at the time of diagnosis					
Race	American Indian/Alaska Native, Asian, Black, Hispanic, White or NA (Not applicable)					
Tumor laterality	Right, left, bilateral					
Oropharynx subsite of origin	The subsite of the tumor within the oropharynx, i.e., base of tongue (BOT), tonsil/soft palate/pharyngeal wall/glossopharyngeal sulcus (GPS)/other (no single subsite of origin could be identified)					
T category	The T category describes the original (primary) tumor, as regard its size and extent, per the American Joint Committee on Gazer (AJCC) and Union for International Cancer Control (UICC) cancer staging system. It could be 'Ti, 'Ti, 'T3, 'T4. https://cancerstaging.org/ references-toolby/gaze/Whati-& Cancer-Staging.apux					
N category	The N category describes whether or not the cancer has reached nearby lymph nodes, per the AJCC and UICC cancer staging system. It can be N0, N1, N2a, N2b, N2c or N3. https://cancerstaging.org/references-tools/Pages/What-is-Cancer-Staging.aspx					
AJCC Stage	AJCC cancer stage. https://cancerstaging.org/references-tools/Pages/What-is-Cancer-Staging.aspx					
Pathological grade	The grade of tumor differentiation. It is described as: I, II, III, IV, I-II, II-III or NA (Not assessable)					
Smoking status at diagnosis	Never, current, or former					
Smoking Pack-Years	An equivalent numerical value of lifetime tobacco exposure. A pack year is defined as twenty cigarettes smoked every day for one year. (NA: Not Assessable) http://smokingpackyears.com/					

# Be FAIR, use the FORCE!



https://www.force11.org/node/6062

#### 3. FAIR FOR MACHINES AS WELL A PEOPLE

In eScience, two clearly separated substrates for knowledge discovery can be distinguished. 1. The actual data, which is as a rule beyond human intellectual capacity to analyse and

2. The 'Explicitome' (everything we already made explicit in text, databases and any other format to date)

- Data should be Findable
- Data should be Accessible
- Data should be Interoperable
- Data should be Re-usable.



# The Future is **FAIR**

# **Comment:** The FAIR Guiding Principles for scientific data management and stewardship

#### Box 2 | The FAIR Guiding Principles

#### To be Findable:

- FIL (meta)data are assigned a globally unique and persistent identifier F2. data are described with rich metadata (defined by R1 below) F3. metadata (derhy and explicitly include the identifier of the data it describes F4. (meta)data are registered or indexed in a searchable resource

- To be Accessible: A1. (meta)data are retrievable by their identifier using a standardized communications protocol A1.1 the protocol is open, free, and universally implementable A1.2 the protocol allows for an authentication and authorization procedure, where necessary A2. metadata are accessible, even when the data are no longer available
- To be Interoperable: 1. (meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation. 12. (meta)data use vocabularies that follow FAIR principles 13. (meta)data include qualified references to other (meta)data

- To be Reusable: R1. meta[data] are richly described with a plurality of accurate and relevant attributes R11. (meta]data are released with a dear and accessible data usage license R1.3. (meta]data meet domain-relevant community standards

# So, we have some work to do...

#### • We need

- Better data aggregation & quality
- More transparent modeling
- Shared datasets
- Standards, standards, standards!
- In short, more FAIR-ness ©

# But our real need is...

• You, the medical physics community...



# Effective data integration is not easy or fast... It's just better!



# We are on the edge of the possible... looking towards the vistas before us!

## Thanks!

Email questions/comments: • <u>cdfuller@mdanderson.org</u>

