# Statistics for Radiomics

Shouhao Zhou, PhD The University of Texas MD Anderson Cancer Center

> Radiomics Certificate, AAPM 2018 Nashville, TN

# Outline

- 1. General statistical modeling concepts and approaches
  - Illustrated in linear regressionObjective: Explanatory vs Predictive
- 2. A statistical predictive feature selection approach for Radiomics
- 3. Machine learning and Statistics
  - Cox regression, cross-validation, LASSO, Bootstrap resampling, ... ...

#### Estimation of Parameters in Linear Model

• The linear model









Objective is to find the estimates for  $\pmb{\beta}_0$  and  $\pmb{\beta}_1$  that minimize the sum of the squared distance from the points to the line.

$$RSS = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i - [\hat{\beta}_0 + \hat{\beta}_1 x_i])^2$$

A good line is one that minimizes the sum of squared differences between the points and the line (Res SS, aka SSE).

SS = "sums of squares"



# Linear Regression – Sources of

Linear Regression - Variation Reg SS Total SS (SSR Due to reg (SST) SST = SSR + SSE $R^2 = SSR/SST$ Res SS



# Overfitting

· Modeling techniques tend to overfit the data.



- Multiple regression:
   Every time you add a variable to the regression, the model's R<sup>2</sup> goes up. Naïve interpretation: every additional predictive variable helps to explain yet more of the target's variance.

  - But that can't be true!
     Error on the dataset used to *fit* the model can be misleading
  - · Doesn't predict future performance.
  - Too much complexity can diminish model's accuracy on future data. • Sometimes called the Bias-Variance Tradeoff.

# Choosing the Optimal Predictive Model

- The model containing all the predictors will always have the smallest RSS and the largest  $R^2,$  since these quantities are related to the training error.
- We wish to choose a model with low test error, not a model with low training error.
  - Training error is usually a poor estimate of test error.
- RSS and R<sup>2</sup> are not suitable for selecting the *best* predictive model among a collection of models with different numbers of predictors.

# Estimating Prediction Error

#### Two approaches:

- Estimate the test error *directly*, by using either a validation set approach or a cross-validation approach.
- Estimate the test error *indirectly*, by making an adjustment to the training error to account for the bias due to overfitting. (e.g. AIC)



# Validation and Overfitting



# Kaplan–Meier curve in training sample (n=129, p=5K)



Subramanian and Simon, JNCI 102, 2010

### K-M curve in test sample



Subramanian and Simon, JNCI 102, 2010

### Cross-Validation (CV, Stone 1974)

+ Partition the full data set D into a prespecified K approximately equal parts  $D_{1'}, D_{2'}, D_{3'}, ..., D_K$ 

 For each k, form a training set T<sub>k</sub> = D - D<sub>k</sub> for model fitting. Compute the validation set error or the CV error for subset D<sub>k</sub> under consideration.

Select the *model* for which the resulting estimated test error is smallest.
 A direct estimate of the test error (prediction error)

- Leave-one-out Cross-validation (LOO-CV) is a special case when K=n
- Can be used in a wider range of model selection tasks

#### 10-fold Cross-validated K-M curve



Simon et al., Briefings in Bioinfo 12, 2011

# **Bias Correction Approaches**

- To assess prediction error, we can use other approaches to correct the bias in estimating the training error:
  - AIC (Akaike information criterion)
  - Mallow's  $C_{\boldsymbol{p}}$  (equivalent to AIC for linear regression)
- Avoid partitioning the data
- These techniques adjust the training error for the model size, and can be used to select among a set of models with different numbers of variables.
- A small value of C<sub>o</sub> and AIC indicates a low error, and thus a better model.

# Mallow's C<sub>p</sub> (Mallow 1973)

For a fitted OLS model containing *d* predictors, the C<sub>p</sub> estimate of test mean squared error:

$$C_p = \frac{1}{n} \left( \text{RSS} + 2d\hat{\sigma}^2 \right)$$

where  $\widehat{\sigma}^2$  is an estimate of the variance of the error  $\epsilon$  associated with each response measurement.

• Here, a penalty is added to the training RSS (Residual SS) in order to adjust for the fact that the training error tends to underestimate the test error.

# Akaike Information Criterion (AIC, Akaike 1973)

· Defined for a large class of models fit by maximum likelihood.

$$AIC = -2\log L + 2 \cdot d$$

where L is the maximized value of the likelihood function for the estimated model.

- In the case of the linear model with Gaussian errors, MLE and OLS are the same things; thus,  $C_{\rm p}$  and AIC are equivalent.
- Stone (1977) proved that in linear regression the error term of both LOO-CV and AIC in estimating the test error are in the same order of o(n<sup>-1</sup>).
   When sample size n -> Inf, the error -> 0.

# LASSO (Tibshirani 1996)

 Recall that the OLS fitting procedure estimates the beta coefficients using the values that minimize:

$$RSS = \sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2$$

 LASSO (least absolute shrinkage and selection operator) is similar to OLS, except that the coefficients are estimated by minimizing a slightly different quantity:

$$\sum_{i=1}^{n} \left( y_i - \beta_0 - \sum_{j=1}^{p} \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^{p} |\beta_j| = \text{RSS} + \lambda \sum_{j=1}^{p} |\beta_j|$$
where  $\lambda \ge 0$  is a component of the determinant of the

# LASSO

- Note that  $\lambda \geq 0$  is a complexity parameter that controls the amount of shrinkage.
- The idea of penalizing by the sum-of-squares of the parameters is also used in neural networks, where it is known as weight decay.
- Lasso uses an  $\ell_1$  penalty, which has the effect of forcing some of the coefficients to be exactly equal to zero when the tuning parameter  $\lambda$  is sufficiently large. Thus, the lasso performs variable/feature selection.
- As  $\lambda$  increases, the standardized ridge regression coefficients shrinks towards zero.
- Thus, when  $\lambda$  is extremely large, then all of the LASSO coefficient estimates are basically zero; this corresponds to the *null model* that contains no predictors.









# Selection of the Tuning Parameter $\boldsymbol{\lambda}$

- Select a grid of potential values; use cross-validation to estimate the error rate on test data (for each value of  $\lambda$ ) and select the value that gives the smallest error rate.
- The model is re-fit using all of the variable observations and the selected value of the tuning parameter  $\lambda.$



#### Survival Analysis

- In many biomedical studies, the primary endpoint is time until an event occurs (e.g. death, recurrence, new symptoms, etc.)
- Data are typically subject to right *censoring* when a study ends before the event occurs.



Often we assume censoring is noninformative, i.e., patients who are censored have the same underlying survival curve after their censoring time as patients who are not censored.

Description of a probability distribution:

1. Distribution function  $F(t) = \Pr(T \le t)$ 2. Density function 3. Survival function 4. Hazard function 5. Cumulative hazard function  $\Lambda(t) = \int_{0}^{t} \frac{\Pr(t \le T \le t + \Delta t \mid T \ge t)}{\Delta t} = \frac{f(t)}{S(t)}$ 

#### Cox PH Regression Model

- Proportional hazards • assumption
- No distributional assumptions on  $\lambda_0(t)$
- Linear effects for • numeric covariates
- Generalized Linear Model

For a single binary covariate, X1:  $\begin{array}{l} X_1 = 0 \implies \\ X_1 = 1 \implies \end{array}$  $\lambda(t; 0) = \lambda_0(t)$  $\lambda(t; 1) = \lambda_0(t) \exp(\beta_1)$ 

• relates hazard function to covariates

$$\begin{split} \lambda(t;\,X) &= \lambda_0(t)\,\,\text{exp}[\,\,\beta_1X_1 + \,... + \,\beta_pX_p] \\ &\bullet \,\lambda_0(t) \text{ is baseline hazard function} \end{split}$$

The ratio of hazard functions,  $HR = \lambda(t; 1)/\lambda(t; 0) = exp(\beta_1)$ Or, more generally  $HR = \lambda(t; X_0 + \Delta X) / \lambda(t; X_0) = \exp(\beta_1 \Delta X)$ 

#### Challenges in Modeling for Radiomics

- Large p small n problem
- Data splitting into training and test sets would reduce the size of the already small training set that is used for the development of the risk prediction model thus increasing problems of instability.

#### Guideline:

- take the information on survival and censored times into consideration.
- binary classification highly depends on the survival threshold used to define the two classes. A slight change of the threshold can lead to very different prediction accuracy and interpretation.
   binary modeling approach can result in loss of efficiency and potential bias in high dimensional settings.
- the choice of threshold affected the predictions.

#### Procedure

#### Pre-screening

- Identify robust and/or correlated radiomics features
- Feature Selection
  - · multivariate Cox proportional hazards regression model with LASSO-penalty

  - CV/GIC to determine the tuning parameters
     Boostrap/CV to determine the most frequently selected features
- Model fitting and Evaluation
  - Comparison between the best clinical model and best clinical+imaging model
- Validation
  - · based on independent test sample

# Other measures for Model Evaluation

- Cross-validated K-M curves
- Harrell's Concordance index (C-index)
   Fraction of all pairs of subjects whose predicted survival times are correctly
  - ordered • Cross-validated C-index for predictive assessment
- Time-dependent ROC curves
- For every specific time t, ROC(t) was plotted as sensitivity(t) versus 1-specificity(t) for all values of the risk score cutoff used to define binary classes
   Optimal cutoff can be determined based Youden's index
- Integrated Brier scores
  - Integrating out over time the score functions which measure the discrepancies between true disease status and predictive risk scores

#### Risk classification based on Statistical Modeling

- 195 Patients, stage III NSCLC w/ definitive XRT
- 11 conventional prognostic factors
- MIM PETedge: Semi-automated delineation
   47 Quantitative Image Features (QIFs) [IBEX]
- Clustering to try to identify multiple risk groups



#### A case when machine learning techniques could be very useful

RTOG 0617 showed no benefit (possible harm) in dose escalation for stage III NSCLC patients
 What if there are sub-groups of patients that would benefit?





#### Impact of Informative Censoring in Survival Analysis

с

- An event is said to be censored if the exact time at which the event occurs is not observed.
- Often we assume censoring is noninformative, i.e., patients who are censored have the same underlying survival curve after their censoring time as patients who are not censored.
- In reality, informative censoring for PFS could arise when patients are censored for initiation of an effective anticancer treatment before the protocol-defined progression.



Campigotto and Weller, JCO 32: 3068-74, 2014

#### Informative Censoring

- Background: programmed cell death protein 1 (PD-1) and programmed deathligand 1 (PD-L1) inhibitors have been increasingly used in cancer therapy.
- Objective: to understand toxicity profile of treatment-related AEs of PD-1/PD-L1 inhibitors
   Approach: to conduct a systematic review
- Approach: to conduct a systematic review and meta-analysis of treatment-related AEs of single-agent PD-1 and PD-L1 inhibitors
- Method: to develop and apply an innovative Bayesian approach for deriving exact inference identical to results using individual-level patient data.
- Result: without correcting the informative censoring, the overall average incidence of all-grade AE would have been 2.65%, an over-estimation of 0.74% or by nearly 40%.



Wang et al. under

review

#### Variance Decomposition

| Factors           |      |     |       |     | Standard Deviations (95% CI) | Rato (95% C)      |
|-------------------|------|-----|-------|-----|------------------------------|-------------------|
| AE                |      |     |       |     |                              |                   |
| Al-grade          |      |     | 0.001 |     | 0.95 [0.90, 1.01]            | 1.50 [1.46. 1.73] |
| Grade 3 and above |      |     | 1-10  |     | 1.17 [1.02. 1.37]            | 1.35 [1.10. 1.67] |
| Drug              |      |     |       |     |                              |                   |
| Až-grade          | - H  |     |       |     | 0.21 [0.05. 0.43]            | 0.35 [0.06, 0.71] |
| Grade 3 and above | - H  | ×   | 4     |     | 0.38 [0.04, 0.79]            | 0.44 (0.04, 0.92) |
| Drug-dose         |      |     |       |     |                              |                   |
| All-grade         | 1.00 |     |       |     | 0.17 [0.08, 0.36]            | 0.29 (0.10, 0.62) |
| Grade 3 and above | -    |     |       |     | 0.17 [0.01, 0.57]            | 0.19[0.01, 0.68]  |
| Cancer            |      |     |       |     |                              |                   |
| All-grade         | -    |     |       |     | 0.04 [0.00, 0.17]            | 0.07 [0.00. 0.29] |
| Grade 3 and above | 1    | •   |       |     | 0.19 [0.01, 0.47]            | 0.22 [0.01, 0.55] |
| Residual          |      |     |       |     |                              |                   |
| All-grade         |      |     |       |     | 0.60 (0.56, 0.64)            |                   |
| Grade 3 and above |      |     |       |     | 0.87 [0.75, 1.00]            |                   |
|                   | -    |     | - 1   |     |                              |                   |
|                   | 0    | 0.5 | 1     | 1.0 |                              |                   |
|                   | 50   |     |       |     |                              |                   |

11



#### Estimation of Longitudinal Biomarker Trajectory

- 635 CML patients received TKI in first/second line treatment.
- 6035 measurements of BCR-ABL expression levels.
- Patients are instructed to schedule follow-up visits approximately at 3, 6, 9, 12, 18, 24 months, and yearly thereafter.
- Real total number of visits varies (Max:23). Real patient visiting times scatter through all the time.



$$\begin{split} f(\mathbf{Y}_{i,\ell} = y_{i,\ell} | \boldsymbol{\mu}, \boldsymbol{\Phi}) &\sim \mathbb{Beta}(\mu_{i,\ell} \phi_{i,\ell}, (1 - \mu_{i,\ell}) \phi_{i,\ell}) \\ & \text{Link}^{0}(\mu_{i,\ell}) = \log i \tau^{-1}(\mu_{i,\ell}) = r_{i,\mu}(t) + \sum_{j} s_{j,\lambda}(X_{j-i,j,l}), \\ & \text{Link}^{0}(\phi_{i,\ell}) = \exp(\phi_{i,\ell}) = r_{i,\mu}(t) + \sum_{j} s_{\phi,\lambda}(X_{\phi,j,\ell}), \\ \text{n specification:} \\ & r_{i,\mu}(t) = \sum_{k} u_{T,\lambda,k,\mu} h_k(t) \\ & \sum_{j} s_{\mu,\lambda}(X_{\mu,\lambda,j}) = \alpha_{i,\mu} + \alpha_{S,\mu} \Re s_k + \sum_{k} u_{A,\lambda,\mu} h_k(\Re s_{\ell}) \end{split}$$

Bayesian Semi-parametric beta regressionEstimation of subject-specific Time-dependent effect

- $$\begin{split} r_{i,\phi}(t) &= \sum_{k} \mathbf{e}_{T,k;\phi} h_{\phi}(t) \\ \Sigma_{f} s_{\phi,j}(X_{\phi,j,i}) &= \alpha_{i,\phi} + \alpha_{S;\phi} \mathbf{Sex}_{i} \end{split}$$
   Patient-Specific Time trend
  - $\mathbf{u}_{T,i,k,\mu} \sim \mathbf{N}(\mathbf{u}_{T,k,\mu}, \sigma_{T,k,\mu}^2)$

Zhou et al. Under review





#### Real-time Prediction





# 2 Cultures in Data Science

#### Machine Learning

- Look for generalizable predictive patterns
- Choose a predictive algorithm by
   Choose a model that relying on its empirical capabilities

Machine learning

- Focus on computational methodology
- Result driven

#### **Traditional Statistics**

- Draw population inferences from a sample
  - incorporates current knowledge of the system
- Focus on math/probability methodology
- Process driven
  - Breiman, *Statist.Sci.* **16**: 199-231, 2001 Blei and Smyth, *PNAS* **114**: 8689-92, 2017 Bzdok et al, *Nature Methods* **15**: 233-34, 2018

#### Glossary

Statistics



| network, graphs           | model                          |
|---------------------------|--------------------------------|
| weights                   | parameters                     |
| learning                  | fitting                        |
| generalization            | test set performance           |
| supervised learning       | regression/classification      |
| unsupervised learning     | density estimation, clustering |
| large grant = \$1,000,000 | large grant= \$50,000          |
|                           |                                |

nice place to have a meeting: nice place to have a meeting: Snowbird, Utah, French Alps Las Vegas in August Source: Rob Tibshirani's personal website

An Era of Data Science

etwor

Models

Fast



Seeing scientific applications turn into methodological advances is always a joy, at least for those of us who care about advancing the science of data, concurrent with advancing science through data.

Model Efficient Sports

earning ence

DeepA

Xiao-Li Meng, Harvard U





# Acknowledgements

 Biostatistics
 Don Berry, PhD
 Xuelin Huang, PhD Xuelin Huang, PhD • Epigenetics Claudio Aldaz, PhD • Emergency Medicine Sai-Ching Yeung, MD PhD FACP • Endocrine Neoplasia Ramona Dadu, MD Ramona Dadu, MD • Experimental Therapeutics Waldemar Priebe, PhD • Imaging Physics Kyle Jones, PhD • Lymphoma Michael Wang, MD • Stem Cell Transfertion Stem Cell Transplantation
Katy Rezvani, MD PhD

#### • Leukemia Kapil Bhalla, MD FACP Jorge Cortes, MD Neurosurgery Frederick Lang, MD FACS Amy Heimberger, MD • Pediatrics Shulin Li, PhD Radiation Physics Laurence Court, PhD Stephen Kry, PhD

• GRA Shuangshuang Fu Nancy Qi Bo Zhao

• 5P30CA016672 • P50CA100632 • P50CA136411 P50CA168505
 R01CA210250 • R01CA214526 • R01CA214749 • R01EB026291 • R21CA216572 • R21CA202104 • HHSN261201500018

 CPRIT Gateway
BMS • MDACC Moonshot