# Machine Learning for Radiomics

Carlos E. Cardenas, Ph.D.

Radiomics Certificate Course – 2018 AAPM Annual Meeting

### Outline

- Introduction
- Data Curation
- Training, Validation, and Test datasets
- Linear and Logistic Regression
- Discriminant Analysis
- Penalized Regression (Ridge and Lasso)
- Support Vector Machines
- Decision Trees and Random Forests
- Clustering Methods
- Final Remarks

THE UNIVERSITY OF TEXAS

erson Cancer

### What is Machine Learning

- "Machine learning is a field of study that gives computers the ability to learn without being explicitly programmed"
  - Arthur Lee Samuel 1959



MDAnderson Cancer Center<sup>®</sup>

# Why Machine Learning?

- Develop systems that can automatically adapt and customize themselves to individual users
- Discover new knowledge from large databases (data mining)
- Automate monotonous tasks (which may require some intelligence)
- Develop systems that are too difficult to hard-code because they require specific detailed skills or knowledge relevant to a specific task
  - Knowledge engineering bottleneck

# Why now?

- Large amounts of clinical data
- Increasing computation power
- Growing progress in available algorithms and theory developed y researchers
- Increasing support from industries and funding agencies

#### Supervised vs Unsupervised Learning



6 Radiomics Certificate Course – 2018 AAPM Annual Meeting

### Supervised vs Unsupervised Learning



THE UNIVERSITY OF TEXAS MDAnderson Cancer Center

#### Supervised Learning



**Goal:** to find specific relationships or structure in the input data that allow us to effectively produce correct output data

#### Unsupervised Learning



**Goal:** to learn the inherent structure of our data without using explicitly-provided labels

# Supervised vs Unsupervised Learning

- Which one should I use?
  - Supervised Learning
    - if you need to train a model to make a prediction--for example, the future value of a continuous variable, such as temperature or a stock price, or a classification—for example, identify makes of cars from webcam video footage.
  - Unsupervised learning
    - if you need to explore your data and want to train a model to find a good internal representation, such as splitting data up into clusters.

UNIVERSITY OF TEXAS

derson Cancer (

#### The two "Trade-offs"

- 1. Prediction Accuracy vs Model Interpretability
  - Global Interpretability
  - Local Interpretability
  - Feature Selection

Interpretability



### The two "Trade-offs"

#### Underfitting **Overfitting** 2. Bias vs Variance High Bias **High Variance** • Bias: error from erroneous assumptions in the learning algorithm Variance: error from Error sensitivity to small Validation Error fluctuations in the training set Expected $\mathrm{E}\left[\left(y-\hat{f}\left(x ight) ight)^{2} ight]=\left(\operatorname{Bias}\left[\hat{f}\left(x ight) ight] ight)^{2}+\mathrm{Var}\left[\hat{f}\left(x ight) ight]+\sigma^{2}$ **Error** Training Error Irreducible where Error $\mathrm{Bias}\left[\hat{f}\left(x ight) ight]=\mathrm{E}\left[\hat{f}\left(x ight)-f(x) ight]$ and Model Complexity $\mathrm{Var}\left[\hat{f}\left(x ight) ight]=\mathrm{E}[\hat{f}\left(x ight)^{2}]-\left(\,\mathrm{E}[\hat{f}\left(x ight)] ight)^{2}$ Source: Bias-Variance Tradeoff in Machine Learning THE UNIVERSITY OF TEXAS MD Anderson Cancer Center<sup>®</sup>

Radiomics Certificate Course – 2018 AAPM Annual Meeting

#### No Free Lunch in Machine Learning



**13** Radiomics Certificate Course – 2018 AAPM Annual Meeting

#### Medical data mining

Linking diseases, drugs, and adverse reactions



14 Radiomics Certificate Course – 2018 AAPM Annual Meeting

Source: Lars Juhl Jensen

MD Anderson Cancer Center<sup>®</sup>

#### Medical data mining

Linking diseases, drugs, and adverse reactions



15 Radiomics Certificate Course – 2018 AAPM Annual Meeting

2012 X Harvard Business Review

#### ANALYTICS

Data Scientist: The Sexiest Job of the 21st Century

#### Cleaning Big Data: Most Time-Consuming, Least Enjoyable Data Science Task, Survey Says



**Gil Press** Contributor (i) Mar 23, 2016, 09:33am • 38,974 views • #BigData





17 Radiomics Certificate Course – 2018 AAPM Annual Meeting

#### 

National Network of Libraries of Medicine

https://nnlm.gov/data/data-curation

#### **Data Curation Tools**

#### (see Data Tools for a more general list of useful tools)

- <u>Open Data Tools: Turning Data into 'Actionable Intelligence'</u> (2013) A comprehensive list of "More than 349 Subject Specific Open Data Tools."
- Information Space: 86 Helpful Tools for the Data Professional PLUS 45 Bonus Tools 
   <sup>®</sup> Very useful
   anthology of tools and resources for data professionals, data dabblers, or data scientists from the iSchool at
   Syracuse.
- Digital Curation Resources outside the DCC № Catalog of tools for data creators and digital curators.
- <u>Digital Curation Glossary</u> № Glossary of data curation and data preservation terminology from the Digital Curation Centre (UK).
- <u>OpenRefine</u> 
   <sup>@</sup> OpenRefine (ex-Google Refine) is a powerful tool for working with messy data, cleaning it, transforming it from one format into another, extending it with web services, and linking it to databases like <u>Freebase</u> 
   <sup>®</sup>.
- <u>ORCID</u> I − An open community-based effort to create and maintain a registry of unique researcher identifiers and a transparent method of linking research activities and outputs to these identifiers.



# Training, Validation, and Test Sets

- Sampling Techniques
  - Simple Random Sampling (SRS)
  - Trial-and-error Methods
  - Systematic Sampling
  - Convenience Sampling
  - CADEX, DUPLEX
  - Stratified Sampling
- Typically in ML/DL we split our data set (of size n) into three subsets:
  - Training
  - Validation
  - Test





THE UNIVERSITY OF TEXAS

derson Cancer (

- Hold-out cross-validation
  - Training, Validation, and Test mutually disjointed datasets
  - Advantage: proportion of subsets are not strictly restricted
  - Train models: Training Set
  - Fine-tune models: Validation Set

#### Algorithm 1 Hold-out cross-validation

- 1. Input: dataset T, performance function *error*, computational models  $L_1, \dots, L_m, m \ge 1$
- 2. Divide T into three disjoint subsets  $T_{tr}$  (training),  $T_v$  (validation), and  $T_t$  (testing).
- 3. For  $j = 1, \dots, m$ :
  - 3.1. Train model  $L_j$  on  $T_{tr}$  and periodically use  $T_v$  to asses the model performance:  $E_v^j = error(L_j(T_v)).$
  - 3.2. Stop training, when a stop-criterion based on  $E_v^j$  is satisfied.

4. For  $j = 1, \dots, m$ , evaluate the performance of the final models on  $T_t$ :  $E_t^j = error(L_j(T_t))$ .

MD Anderson Cancer Center<sup>®</sup>

THE UNIVERSITY OF TEXAS

- k-fold cross-validation
  - More test  $\rightarrow$  stable estimate of the model error
  - Useful when not enough data is available

#### Algorithm 2 K-fold cross-validation

- 1. Input: dataset T, number of folds k, performance function error, computational models  $L_1,\cdots,L_m,m\geq 1$
- 2. Divide T into k disjoint subsets  $T_1, \dots, T_k$  of the same size.

3. For 
$$i = 1, \dots, k$$
:  
 $T_v \leftarrow T_i, T_{tr} \leftarrow \{T \setminus T_i\}.$ 

3.1. For  $j = 1, \dots, m$ :

Train model  $L_j$  on  $T_{tr}$  and periodically use  $T_v$  to asses the model performance:

$$E_v^j(i) = error(L_j(T_v)).$$

Stop training, when a stop-criterion based on  $E_v^j(i)$  is satisfied.

4. For  $j = 1, \dots, m$ , evaluate the performance of the models by:  $E_v^j = \frac{1}{k} \cdot \sum_{i=1}^k E_v^j(i)$ .

Resource: <u>Data Splitting</u> by Z. Reitermanova

MD Anderson Concer Center<sup>®</sup>

THE UNIVERSITY OF TEXAS

• Combining hold-out and k-fold CV



22 **Radiomics Certificate Course – 2018 AAPM Annual Meeting** 

• Combining hold-out and k-fold CV



MDAnderson Cancer Center<sup>®</sup>

• Combining hold-out and k-fold CV



# Machine Learning Models

25 Radiomics Certificate Course – 2018 AAPM Annual Meeting

#### "Some food for thought"

• George Box (1919 – 2013)



#### "Essentially, all models are wrong, but some are useful"

#### 26 Radiomics Certificate Course – 2018 AAPM Annual Meeting

#### Linear Regression



#### Linear Regression

- Assessing the accuracy of the model
  - R-squared or fraction of variance explained

$$R^{2} = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$
$$\text{TSS} = \sum_{i=1}^{n} (y_{i} - \bar{y})^{2}$$

$$RSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

• Residual Sum of Squares

$$\mathrm{RSE} = \sqrt{\frac{1}{n-2}}\mathrm{RSS}$$

#### Linear Regression

- Assumptions
  - Linear relationship
    - Residual plots
  - Multivariate normality (predictors are normally distributed)
    - Goodness of fit test (Kolmogorov-Smirnov test)
  - No or little multicollinearity
    - Correlation matrix (Pearson's), Tolerance, and variance inflation factor (VIF)

THE UNIVERSITY OF TEXAS

erson Cancer

- No auto-correlation
  - Time-series data
- Homoscedasticity (residuals are equal across the regression line)
  - Scatter plots

- Less assumptions than linear regression
- Some still apply
  - Observations to be independent of each other
  - Little or no multicollinearity
  - Linearity of independent variables and log odds
  - Larger sample size is useful
- The only "real" limitation on logistic regression is that the outcome must be discrete

THE UNIVERSITY OF TEXAS

derson Cancer (`enter®

• Discrete Outcomes (Classification Problem)

$$p(X) = \frac{e^{\beta_0 + \beta_1 X}}{1 + e^{\beta_0 + \beta_1 X}} \quad \text{Logistic Function}$$

#### **Estimation of Parameters:**

Maximum likelihood:

$$\ell(\beta_0, \beta_1) = \prod_{i:y_i=1} p(x_i) \prod_{i':y_{i'}=0} (1 - p(x_{i'}))$$

Likelihood Function

• Discrete Outcomes (Classification Problem)



32 Radiomics Certificate Course – 2018 AAPM Annual Meeting

- Limitations
  - Unstable with well separated classes
  - Unstable with few examples



#### Discriminant Analysis

• Classification algorithm that estimates Bayesian classification



$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

Multi-variable solution

$$\hat{\Sigma}_k := \frac{1}{\hat{n}_k} \sum_{i:y_i=k} (x_i - \hat{\mu}_k) (x_i - \hat{\mu}_k)^T$$

#### Discriminant Analysis

• Discriminant analysis predicts as follows:

 $\hat{Y}|X = x := \operatorname{argmax}_k \pi_k p_k(x) = \operatorname{argmax}_k \delta_k(x)$ 



#### Discriminant Analysis

- Key Assumptions
  - Each class density is multivariate Gaussian

$$X|Y_j \sim N(\mu_j, \Sigma_j), \quad j = 0, 1$$

• Equal covariance

$$\Sigma_j = \Sigma, \quad j = 0, 1$$

THE UNIVERSITY OF TEXAS

Anderson Cancer Center<sup>®</sup>

• No outliers

# Ridge Regression

- Like least squares linear regression but *shrinks* the estimated coefficients towards zero
- Very useful when multicollinearity (near-linear relationships among the independent variables) occurs
- Given a response vector  $y \in \mathbb{R}^n$  and a predictor matrix  $X \in \mathbb{R}^{n \times p}$ , the ridge regression coefficients are defined as

$$\hat{\beta}^{\text{ridge}} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2$$
$$= \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\|y - X\beta\|_2^2}_{\operatorname{Loss}} + \lambda \underbrace{\|\beta\|_2^2}_{\operatorname{Penalty}} \longrightarrow \mathsf{L}_2 \operatorname{norm}$$

THE UNIVERSITY OF TEXAS

MD Anderson Concer Center<sup>®</sup>

• Where  $\lambda \geq 0$  is a tuning parameter that controls the strength of the penalty term

# Ridge Regression

- Cannot perform variable selection
  - Coefficients are reduced close to zero, but not zero (unless λ = ∞, where all coefficients are zero)
  - $\rightarrow$  Low interpretability



# Ridge Regression

• Variable standardization is the first step when using ridge regression!

THE UNIVERSITY OF TEXAS

Anderson <del>Cancer</del> Center®

- •Assumptions
  - Linearity
  - Constant variance (no outliers)
  - Independence

- The Lasso combines some of the shrinking advantages of ridge regression with variable selection
- It is very competitive with the ridge regression in regards to prediction error
- The only difference between the two is that ridge regression uses  $\ell_2$  norm penalty where the lasso uses the  $\ell_1$  norm penalty
- While the  $\ell_1$  and  $\ell_2$  norm look very similar, the ridge and lasso solutions behave very differently

• The Lasso (Least Absolute Selection and Shrinkage Operator) is defined

$$\hat{\beta}^{\text{lasso}} = \underset{\beta \in \mathbb{R}^{p}}{\operatorname{argmin}} \frac{\|y - X\beta\|_{2}^{2}}{\|y - X\beta\|_{2}^{2}} + \lambda \sum_{j=1}^{p} |\beta_{j}|$$
$$= \underset{\beta \in \mathbb{R}^{p}}{\operatorname{argmin}} \underbrace{\|y - X\beta\|_{2}^{2}}_{\operatorname{Loss}} + \lambda \underbrace{\|\beta\|_{1}}_{\operatorname{Penalty}} \longrightarrow L_{1} \operatorname{norm}$$

Anderson Concer (Center®

- Again, we have a tuning parameter  $\boldsymbol{\lambda}$  that controls the amount of regularization
- As usual, assume  $X^{n \times p}$  is standardize and y is centered

- Often, we believe that many of the  $\beta_j$  's should be 0
- Therefore, we would like to have a set of **sparse solutions**
- Large enough  $\lambda$  will set some coefficients exactly equal to 0!
  - So the LASSO will perform variable selection for us!



THE UNIVERSITY OF TEXAS

MD Anderson Concer

*enter* 



- Limitation
  - Cases where p >> n the lasso selects at most n variables before it saturates
    - Convex optimization problem
    - Curse of dimensionality
- Solution
  - Elastic Nets (Ridge + Lasso)
  - Zou, Hui; Hastie, Trevor (2005). "Regularization and Variable Selection via the Elastic Net". *Journal of the Royal Statistical Society*. Series B. Wiley. **67**(2): 301–20. JSTOR 3647580

THE UNIVERSITY OF TEXAS

Anderson Concer (Center

# Classification w/ Ridge and Lasso

• Add a penalty to the logistic function

$$l_p(\beta_0; \beta; \lambda) = -l(\beta_0; \beta) + \lambda J(\beta)$$

THE UNIVERSITY OF TEXAS

Anderson Concer Center

- Where
  - $l(eta_0;eta)$  denotes the unrestricted log-likelihood function
  - $\lambda$  is the regularization parameter controlling the amount of shrinkage
  - $J(\beta)$  is a penalty function on the coefficient parameter  $\beta$ 
    - Either the lasso or ridge penalty functions

- Goal: find an optimal hyperplane
- Support Vectors
  - Data points that lie closest to the decision surface (or hyperplane)
  - Most difficult data points to classify
  - They have a direct relationship on the optimal location of the hyperplane



- In general, there are lots of possible solutions
- Support Vector Machine (SVM) finds an <u>optimal</u> solution

- SVMs <u>maximize</u> the margin around the separating hyperplane
- The decision function is fully specified by a subset of training samples
  - Support Vectors
- Real-world data (non-separable)
  - Soft margin classifier
    - OK to misclassify a few training observations in order to do a better job in classifying the remaining observations



$$egin{aligned} ext{maximize} & f(c_1 \dots c_n) = \sum_{i=1}^n c_i - rac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i (arphi(ec{x}_i) \cdot arphi(ec{x}_j)) y_j c_j \ & = \sum_{i=1}^n c_i - rac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i c_i \overline{k(ec{x}_i, ec{x}_j)} y_j c_j \ & ext{Kernel} \ & ext{subject to} & \sum_{i=1}^n c_i y_i = 0, ext{ and } 0 \leq c_i \leq rac{1}{2n\lambda} ext{ for all } i. \end{aligned}$$

- The Kernel Trick
  - Linear:  $k(\overrightarrow{x_i}, \overrightarrow{x_j}) = (\overrightarrow{x_i} \cdot \overrightarrow{x_j})$
  - Polynomial:  $k(\overrightarrow{x_i}, \overrightarrow{x_j}) = (\overrightarrow{x_i} \cdot \overrightarrow{x_j})^p$
  - Gaussian radial basis function:  $k(\vec{x_i}, \vec{x_j}) = e^{-\gamma \|\vec{x_i} \vec{x_j}\|^2}$  for  $\gamma > 0$
  - Hyperbolic tangent:  $k(\vec{x_i}, \vec{x_j}) = \tanh(\kappa \vec{x_i} \cdot \vec{x_j} + c)$  for some  $\kappa > 0$  and c < 0

- Kernel Trick
  - Useful when the decision function is not linear



#### What to do when your SVM is overfitting?

- Decrease the cost constraint (C) of the regularization term in the Lagrange formulation
- Use less expressive kernel (maybe a smaller degree polynomial)
- As with any learner, collecting more training data usually provides an improvement in accuracy

UNIVERSITY OF TEXAS

derson Cancer

- CART (Classification and Regression Trees, Breiman et al 1984)
- A decision tree is drawn upside down with its root at the top
  - Node
  - Branches
  - Leaf (decisions)



How does the algorithm decide which features are more important (top of the tree) and where to create a split?

- Recursive binary splitting
  - All the features are considered and different split points are tried and tested using a cost function

derson Cancer ( 'er

- The split with the lowest cost is selected
- "Greedy algorithm"
- Cost functions
  - Regression: Mean Square Error
  - Classification: Gini, Entropy, etc.

#### Improving prediction accuracy

- Assigning a maximum depth
  - Length of the longest path from the root to leaf
- Leaf size
  - Set a minimum number of training inputs to use on each leaf
- Pruning
  - Removal of branches that make use of features having low importance

Reduces complexity of the tree  $\rightarrow$  increased predictive power by reducing overfitting

 In theory, the depth of the tree is limited by the number of training examples and extremely deep

nderson Cancer (Center®

- Advantages
  - Simple to understand, interpret, visualize
  - Handles both numerical and categorical data
  - Nonlinear relationships between parameters do not affect tree performance
- Disadvantages
  - Decision tree learners can create over-complex trees  $\rightarrow$  overfitting
  - Decision tree learners create biased trees if some classes dominate
    - Balance the data set prior to training
  - Greedy algorithms cannot guarantee to return the globally optimal decision tree → mitigated by ensemble learning (multiple trees)

derson Concer

• Adaboost, Bagging, Random Forest, etc.



https://www.eustafor.eu/



• Adds additional randomness to the model



- 1. For each tree: randomly selects a subset of training data (~66%)
- 2. At each node: randomly selects a subset of predictor variables ( $\sim \sqrt{p}$ )

THE UNIVERSITY OF TEXAS

MD Anderson Cancer Center<sup>®</sup>

- Out-of-Bag (OOB) Error
  - Estimates the prediction error of random forests (and other ensemble learners) by using only the trees that did not have a training sample  $x_i$  in their bootstrap sample
- Feature importance
  - Very easy to calculate
  - Can be used to remove low importance features

THE UNIVERSITY OF TEXAS

nderson Cancer

- Advantages
  - Both regression and classification
  - Easy to use
  - Number of hyperparameters is not too high and they are easy to understand
  - More trees  $\rightarrow$  better predictions (until you reach a plateau)
- Disadvantages
  - A lot of trees  $\rightarrow$  slow and ineffective for real-time predictions

derson Cancer

• Loss of interpretability

# Clustering Methods

- Unsupervised learning technique
- Helps identify homogenous subgroups or clusters in a data set

derson Cancer

- Two clustering approaches:
- K-Means Clustering
- Hierarchical Clustering

#### **K-Means** Clustering

• How it works

Algorithm 10.1 K-Means Clustering

- 1. Randomly assign a number, from 1 to K, to each of the observations. These serve as initial cluster assignments for the observations.
- 2. Iterate until the cluster assignments stop changing:
  - (a) For each of the K clusters, compute the cluster *centroid*. The kth cluster centroid is the vector of the p feature means for the observations in the kth cluster.
  - (b) Assign each observation to the cluster whose centroid is closest (where *closest* is defined using Euclidean distance).

Source: An Introduction to Statistical Learning, Witten et al, 2013





Source: An Introduction to Statistical Learning, Witten et al, 2013

### **K-Means** Clustering

- Limitations
  - Random initialization & local optima. R



Source: An Introduction to Statistical Learning, Witten et al, 2013

### **Hierarchical Clustering**

• Creates a dendrogram

#### Algorithm 10.2 Hierarchical Clustering

- 1. Begin with n observations and a measure (such as Euclidean distance) of all the  $\binom{n}{2} = n(n-1)/2$  pairwise dissimilarities. Treat each observation as its own cluster.
- 2. For  $i = n, n 1, \dots, 2$ :
  - (a) Examine all pairwise inter-cluster dissimilarities among the *i* clusters and identify the pair of clusters that are least dissimilar (that is, most similar). Fuse these two clusters. The dissimilarity between these two clusters indicates the height in the dendrogram at which the fusion should be placed.
  - (b) Compute the new pairwise inter-cluster dissimilarities among the i-1 remaining clusters.

Source: An Introduction to Statistical Learning, Witten et al, 2013



THE UNIVERSITY OF TEXAS MDAnderson Cancer Center®

# Clustering Methods

- •Things to remember when using clustering algorithm:
  - Standardizing variables so that all are on the same scale. It is important when calculating distances
  - Treat data for outliers before forming clusters as it can influence the distance between the data points.

UNIVERSITY OF TEXAS

nderson Cancer

### Honorable Mention List

- •Elastic Nets
- Principal Component Analysis

THE UNIVERSITY OF TEXAS

Anderson Cancer

- •Other Ensemble Methods
  - Bagging and Boosting
- •Nearest Neighbor
- •Bayesian Networks

#### Final Remarks

- Proper splitting of datasets leads to better generalization
  - No peaking at the test set!
- No one algorithm works best for every problem
  - "No free lunch in Machine Learning"
- Think about the inputs in your model
  - Intuition and knowledge about the data can prevent head-aches
- Develop a good understanding about mathematical principals behind your algorithm of choice!

on Concer

• Identify the strengths and weaknesses

# Thank you!

Carlos E. Cardenas, PhD cecardenas@mdanderson.org

66 Radiomics Certificate Course – 2018 AAPM Annual Meeting