

Methods for Analytical Validation

William Hsu, PhD

Medical & Imaging Informatics group

Associate Professor of Radiological
Sciences, Bioinformatics, and

Bioengineering

whsu@mednet.ucla.edu

Learning Objectives

- Describe the process of analytical validation
- Explain basic concepts and metrics
- Develop a study for evaluating and comparing algorithms
- Assess the results of an evaluation

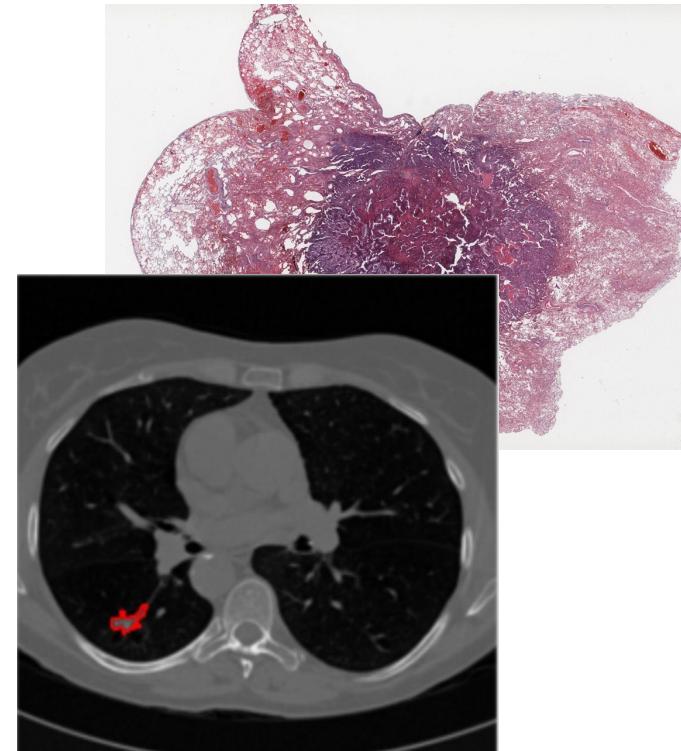
Headings in “green” were included in the printed course notes

Headings in “blue” are additional slides that are available in the updated course notes (<https://uclawillhsu.github.io/spie2023mi/>).

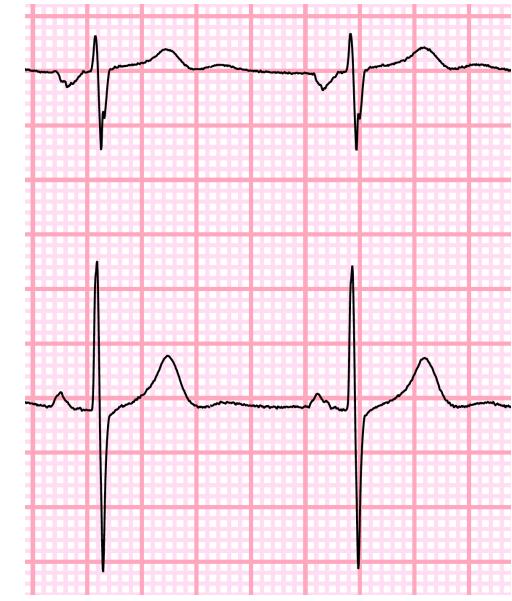
Learning from varied data types...

Patient ID	Chest Pain	Myocardial Infarct
1234	YES	YES
1235	NO	NO
1236	YES	NO
1237	YES	YES
1238	NO	NO

Textual / numerical

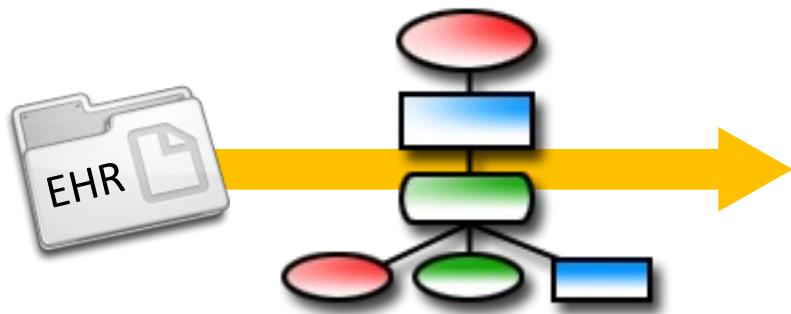


Imaging

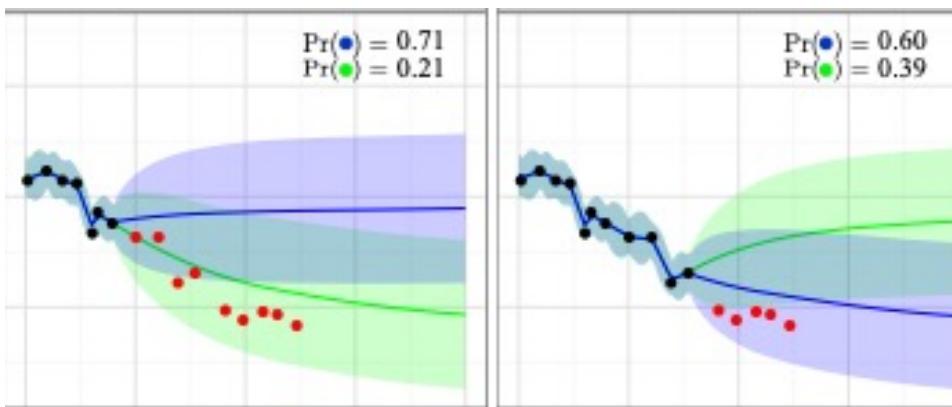


Waveform

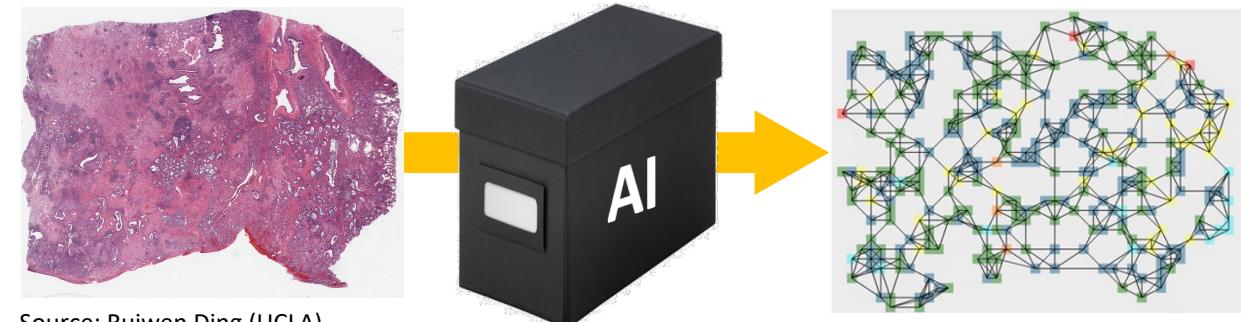
... to provide actionable insights



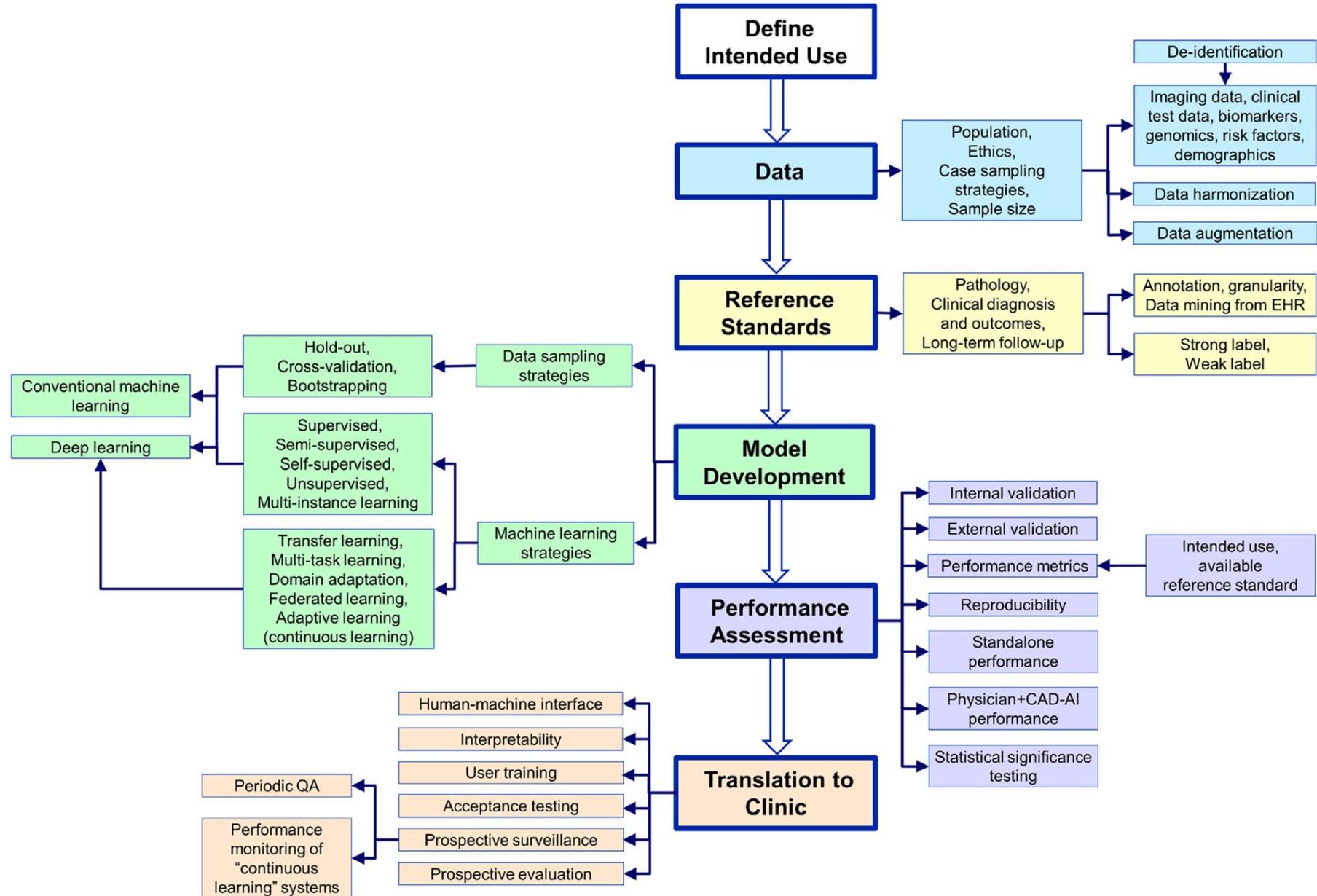
Diagnosis?
Clinical deterioration?



Source: <https://arxiv.org/pdf/1601.04674.pdf>

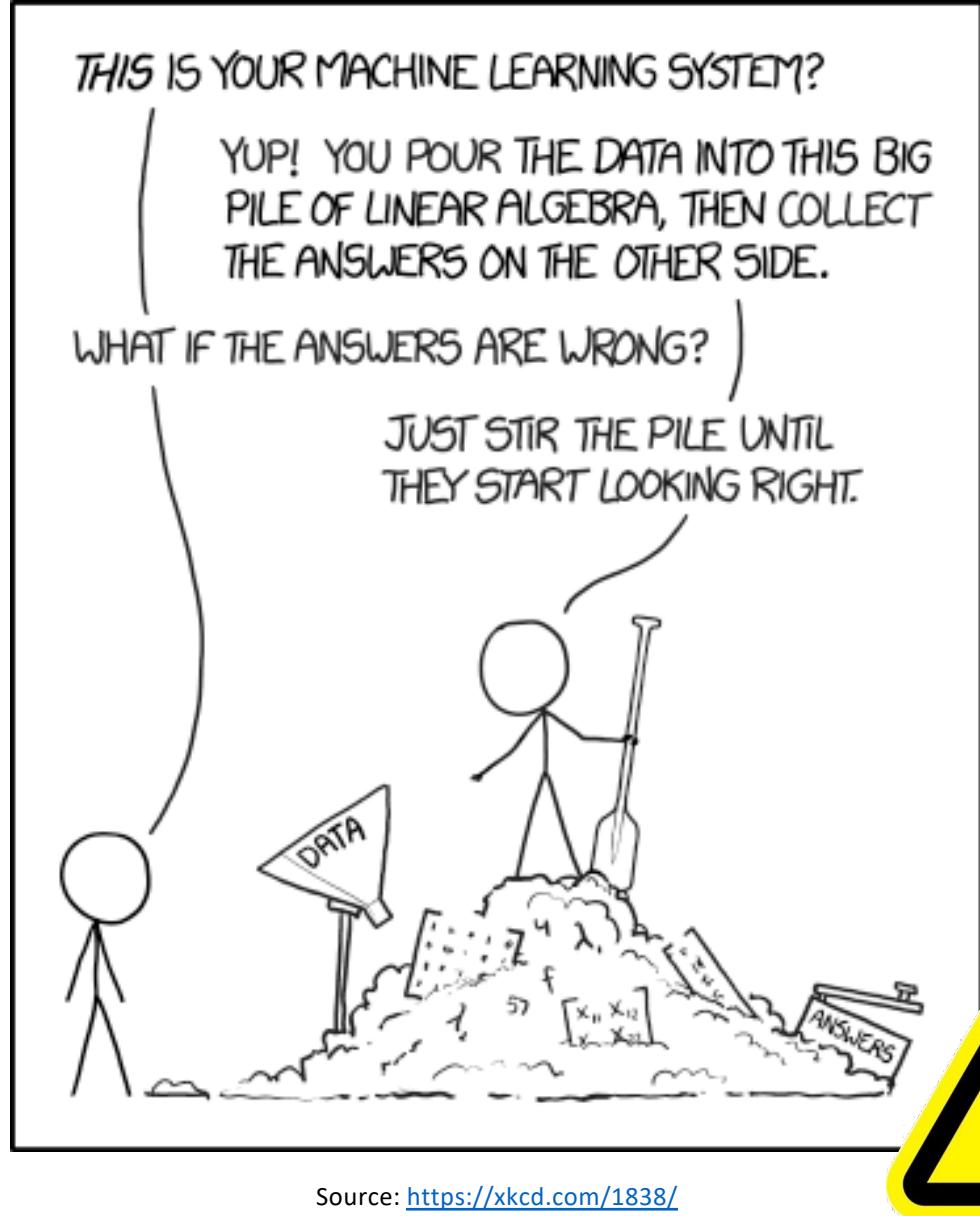


Source: Ruiwen Ding (UCLA)



Driving questions

1. How do we assess whether a model performs as intended?
2. How do we select between multiple models that accomplish the same task?
3. Will the model advance medical decision making?



Source: <https://xkcd.com/1838/>



Analytical Validation

- Does the learning algorithm correctly process input data to generate accurate, reliable, and precise output data?

Clinical Evaluation		
Valid Clinical Association	Analytical Validation	Clinical Validation
Is there a valid clinical association between your SaMD output and your SaMD's targeted clinical condition?	Does your SaMD correctly process input data to generate accurate, reliable, and precise output data?	Does use of your SaMD's accurate, reliable, and precise output data achieve your intended purpose in your target population in the context of clinical care?

Figure 1 - Clinical Evaluation Process

SaMD: Software as a Medical Device

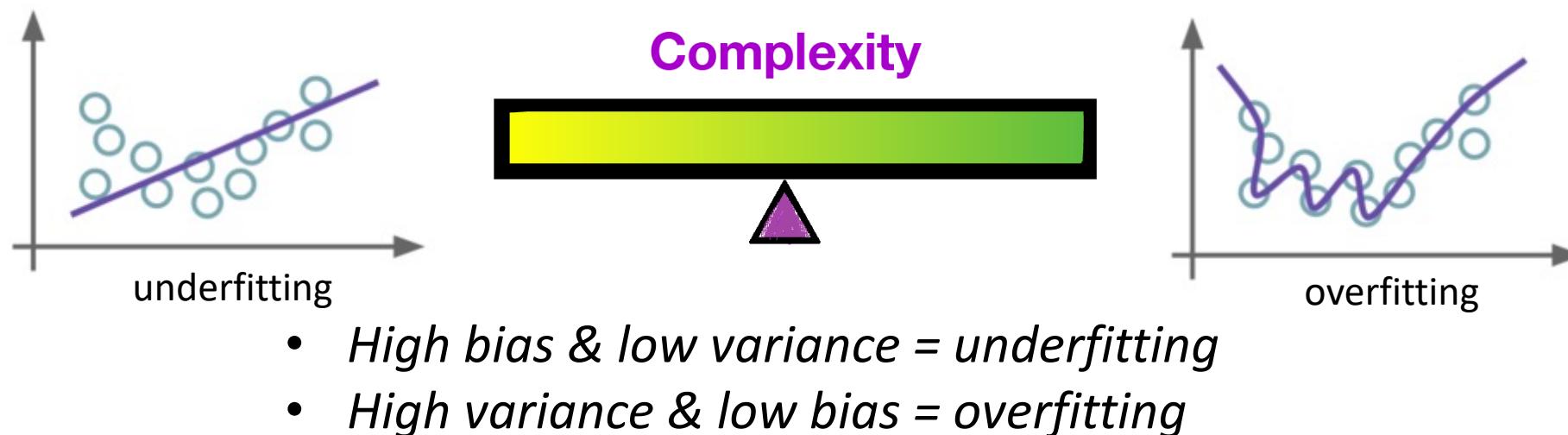
Source: <https://www.fda.gov/files/medical%20devices/published/Software-as-a-Medical-Device-%28SaMD%29--Clinical-Evaluation---Guidance-for-Industry-and-Food-and-Drug-Administration-Staff.pdf>

Basic concepts: Prediction error

- **Prediction error** can refer to one of two things:
 - In **regression analysis**, it is a measure of how well the model predicts the response variable
 - E.g., root mean squared error
 - In **classification**, it is a measure of how well samples are classified to the correct category
 - E.g., accuracy, precision, recall

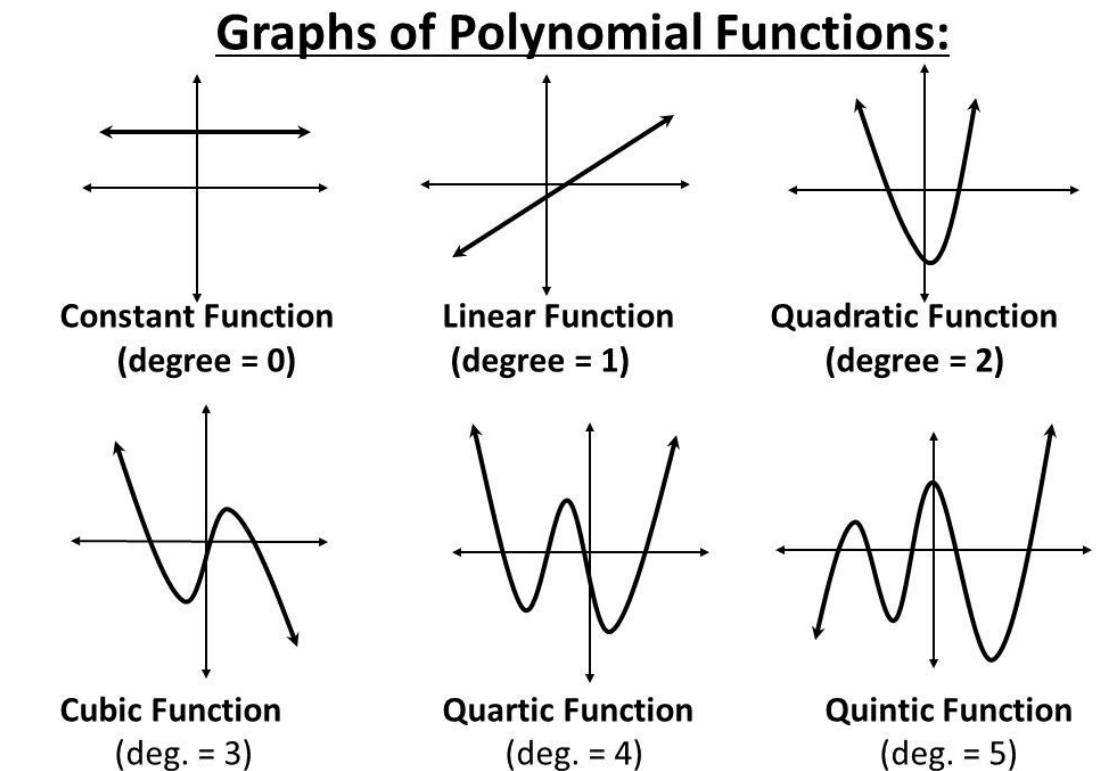
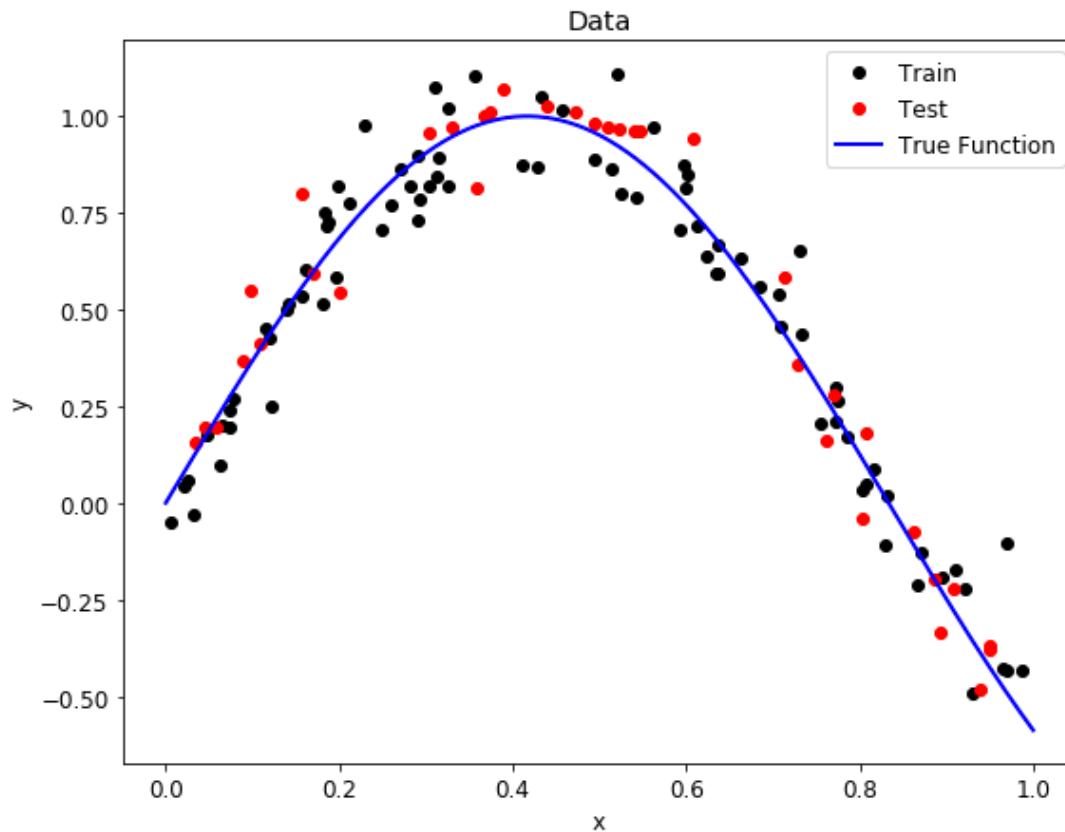
Basic concepts: Bias-variance tradeoff

- **Bias** is the difference between the average prediction of our model and the correct value we are trying to predict.
- **Variance** is the model's sensitivity to small fluctuations in the training set. High variance can cause an algorithm to model the random noise in the training data rather than the intended outputs.



Regression example

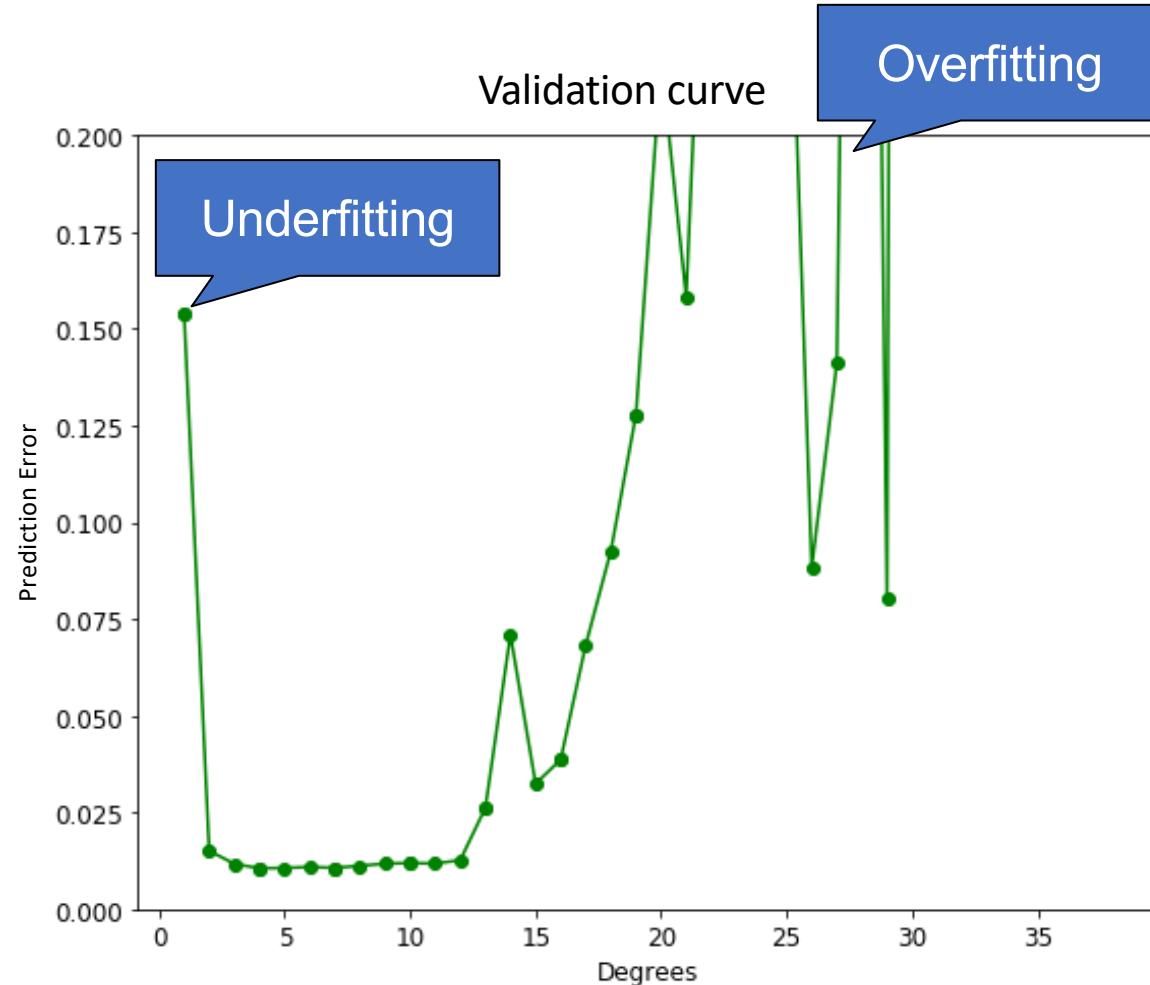
Evaluate polynomial function that best fits the data



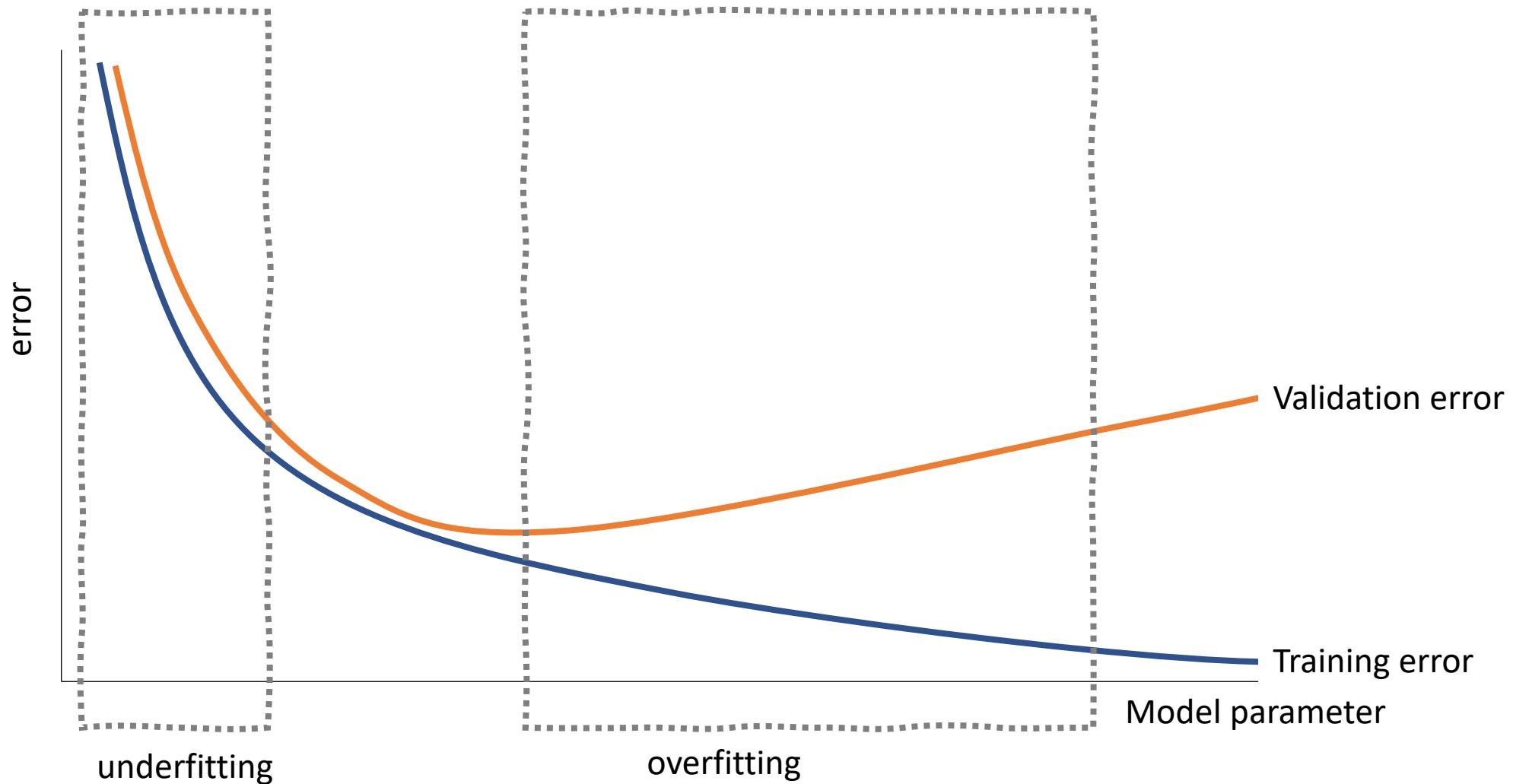
Regression example

Validation results
ordered by
prediction error

Degrees	MSE
4	0.010549
5	0.010637
7	0.010665
6	0.010887
8	0.011182
3	0.011695
9	0.011757
11	0.011769
10	0.011902
12	0.012642



Training and validation curves



Confusion matrix

2-class problems

		Predicted Class	
		Positive	Negative
Real Class	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{FP} + \text{FN} + \text{TN}}$$

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

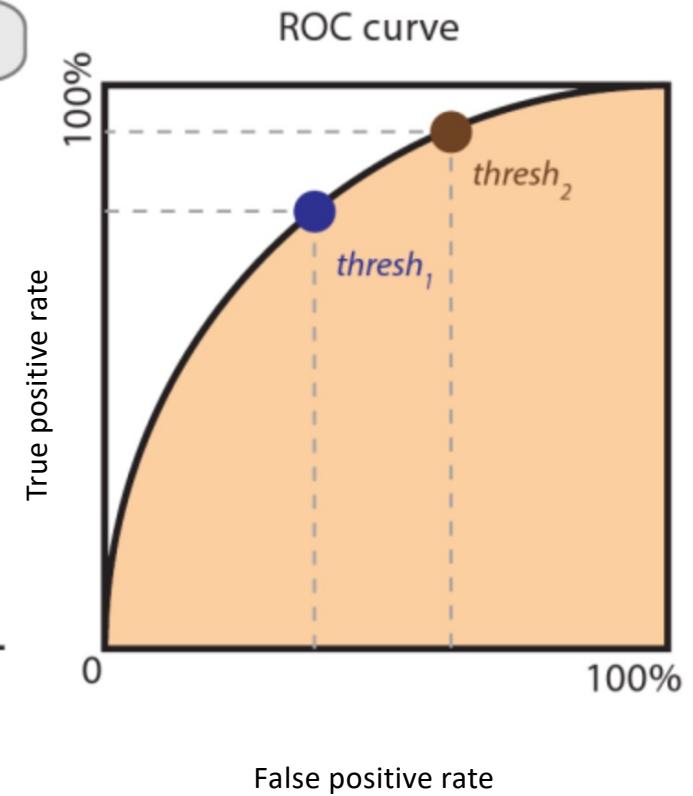
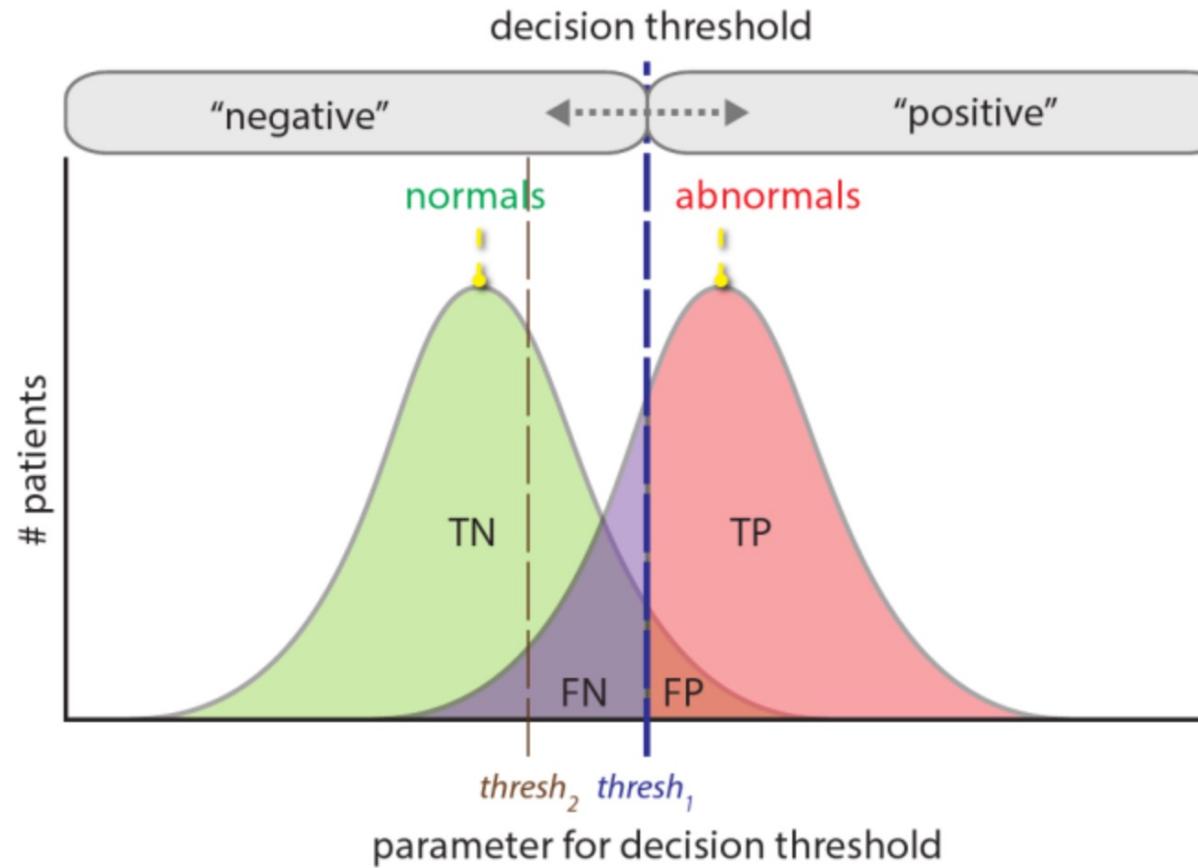
$$\text{true positive rate (recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{false positive rate} = \frac{\text{FP}}{\text{TN} + \text{FP}} \\ (1 - \text{specificity})$$

Receiver Operator Characteristic Curve

$$\text{true positive rate (recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{false positive rate} = \frac{\text{FP}}{\text{TN} + \text{FP}} \quad (1 - \text{specificity})$$



Source: Sun, Shi, and Mandell. Core Radiology 2nd ed. Chapter 15. <https://doi.org/10.1017/9781108966450>


```

.lroc

Logistic model for var1

number of observations =          10
area under ROC curve = 0.7400

.lstat

Logistic model for var1

      _____ True _____
      |           | D   | ~D |
Classified |           |-----|-----|
      +       |       3   |   1  |       4
      -       |       2   |   4  |       6
      |           |-----|-----|
      Total    |       5   |   5  |      10

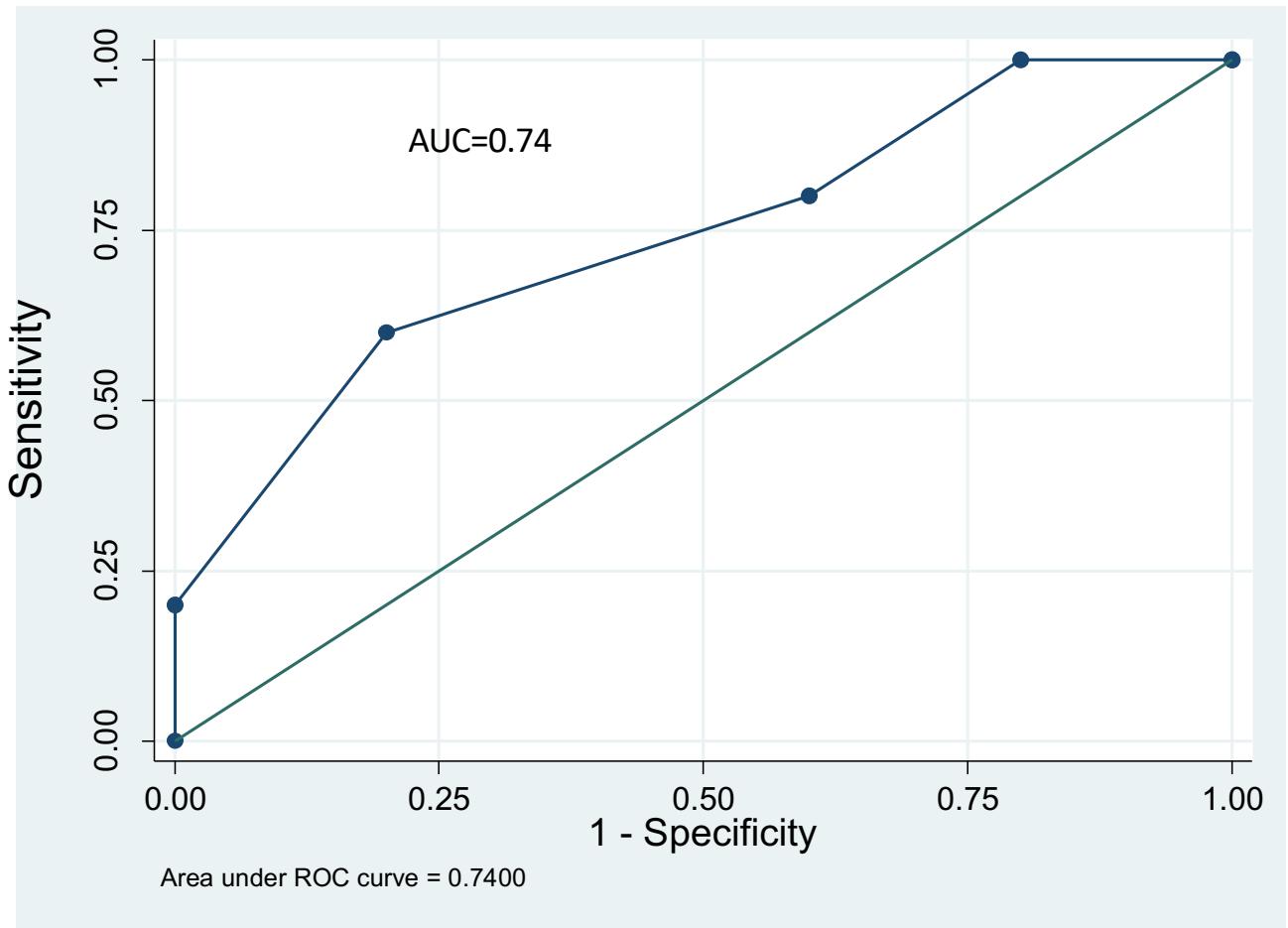
Classified + if predicted Pr(D) >= .5
True D defined as var1 != 0

Sensitivity          Pr( + | D)  60.00%
Specificity          Pr( - | ~D) 80.00%
Positive predictive value  Pr( D | +) 75.00%
Negative predictive value  Pr(~D | -) 66.67%

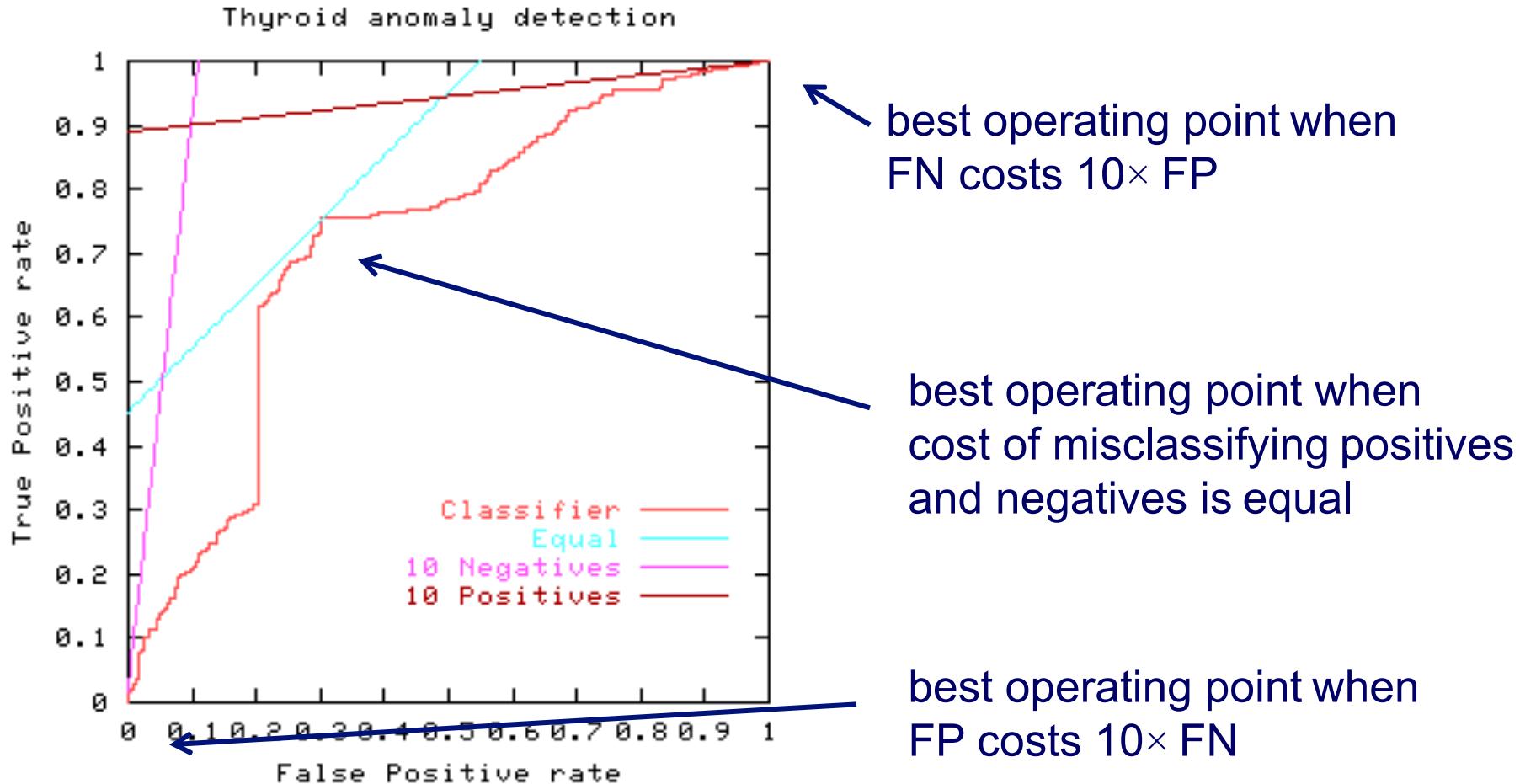
False + rate for true ~D  Pr( + | ~D) 20.00%
False - rate for true D  Pr( - | D) 40.00%
False + rate for classified +  Pr(~D | +) 25.00%
False - rate for classified -  Pr( D | -) 33.33%

Correctly classified 70.00%

```

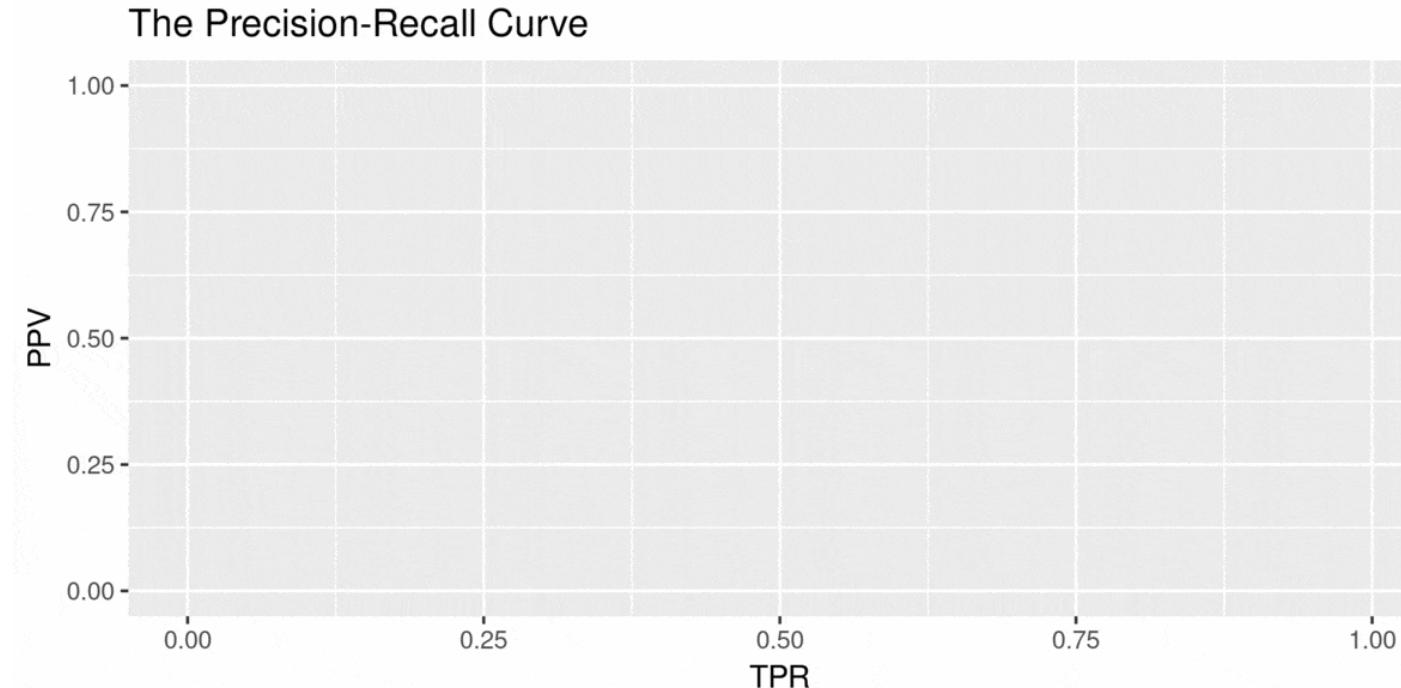


ROC curves and misclassification costs



Source: David Page (University of Wisconsin)

Precision-recall curves

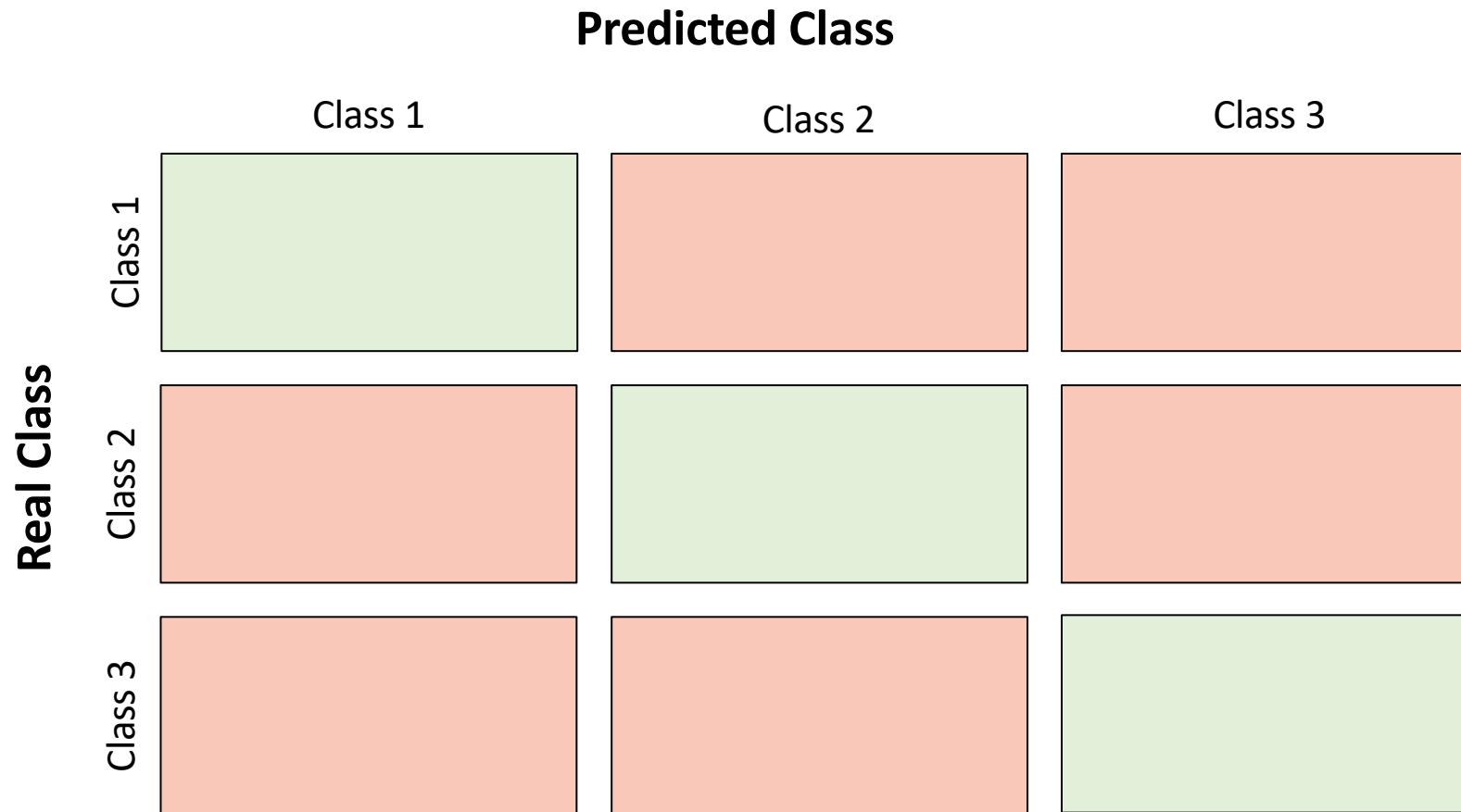


Plots the precision vs. recall as a threshold on the confidence of an instance being positive is varied

<i>Estimate</i>	-3.0	-2.5	-1.0	-1.0	-0.5	-0.5	0.5	1.0	2.0	2.0	3.5
<i>Class</i>	-1	-1	-1	-1	-1	+1	-1	+1	-1	+1	+1
<i>Prediction</i>	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA
<i>Status</i>	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA

Source: <https://www.datascienceblog.net/post/machine-learning/interpreting-roc-curves-auc/>

Multi-class scenario



Example

		True/Actual		
		Cat (😺)	Fish (🐠)	Hen (🐓)
Predicted	Cat (😺)	4	6	3
	Fish (🐠)	1	2	0
	Hen (🐓)	1	2	6

<https://towardsdatascience.com/multi-class-metrics-made-simple-part-i-precision-and-recall-9250280bddc2>

$$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{true positive rate (recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Precision of 'cat' classifier
 $= 4/(4+6+3) = 30.8\%$

Recall of 'cat' classifier
 $= 4/(4+1+1) = 66.7\%$

Overall accuracy
 $= (4+2+6)/25 = 48\%$

Multi-class scenario

- **No averaging:** Return a measure corresponding to each class.

		True/Actual		
		Cat (😺)	Fish (🐠)	Hen (🐓)
Predicted	Cat (😺)	4	6	3
	Fish (🐠)	1	2	0
	Hen (🐓)	1	2	6

<https://towardsdatascience.com/multi-class-metrics-made-simple-part-i-precision-and-recall-9250280bddc2>

- Precision of 'cat' classifier = $4/(4+6+3) = 30.8\%$
- Precision of 'fish' classifier = $2/(1+2) = 66.7\%$
- Precision of 'hen' classifier = $6/(1+2+6) = 66.7\%$

Multi-class scenario

- **Macro averaging:** Calculate the measure independently for each class then find their unweighted mean.

		True/Actual		
		Cat (😺)	Fish (🐠)	Hen (🐓)
Predicted	Cat (😺)	4	6	3
	Fish (🐠)	1	2	0
	Hen (🐓)	1	2	6

<https://towardsdatascience.com/multi-class-metrics-made-simple-part-i-precision-and-recall-9250280bddc2>

- Macro-precision of classifier is $(31\% + 67\% + 67\%) / 3 = 54.7\%$

Multi-class scenario

- **Micro averaging:** Aggregate contributions of all classes to compute the average measure by counting the total number of times each class was correctly predicted and incorrectly predicted.

		True/Actual		
		Cat (😺)	Fish (🐠)	Hen (🐓)
Predicted	Cat (😺)	4	6	3
	Fish (🐠)	1	2	0
	Hen (🐓)	1	2	6

<https://towardsdatascience.com/multi-class-metrics-made-simple-part-i-precision-and-recall-9250280bddc2>

- $TP=4+2+6=12$
- $FP=6+3+1+0+1+2=13$
- Micro-precision of classifier is $12/(12+13) = 48\%$

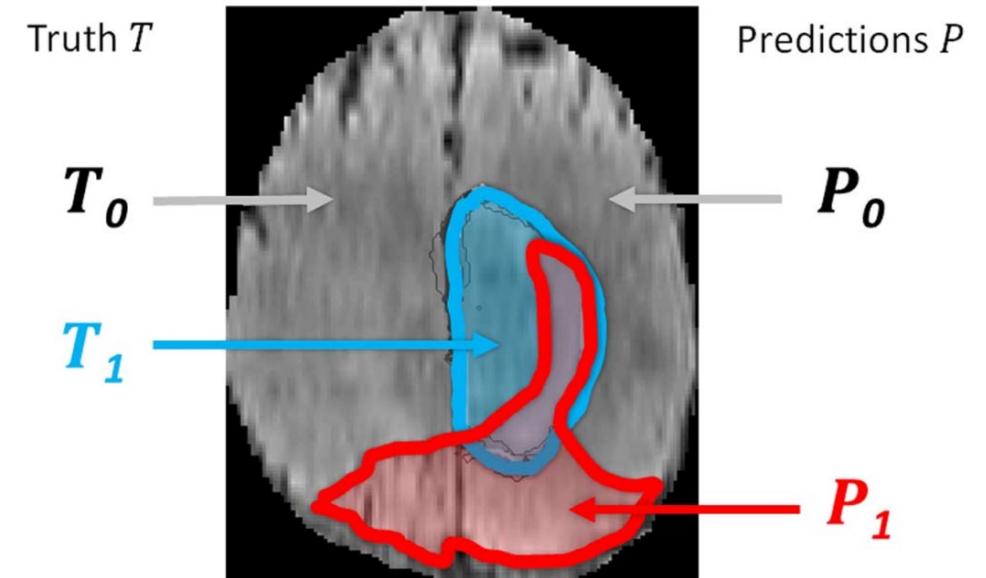
Notes about measures

- Use precision and recall to focus on small positive class
 - If correct detection of negatives examples is less important to the problem, use precision and recall
- Use ROC when both classes detection is equally important
- Use ROC when the positives are the majority
 - Precision and recall reflects the ability of prediction of the positive class and not the negative class which will naturally be harder to detect due to the smaller number of samples.
- Micro-average is preferable over macro-average in situations where a class imbalance exists

Other potential measures of interest

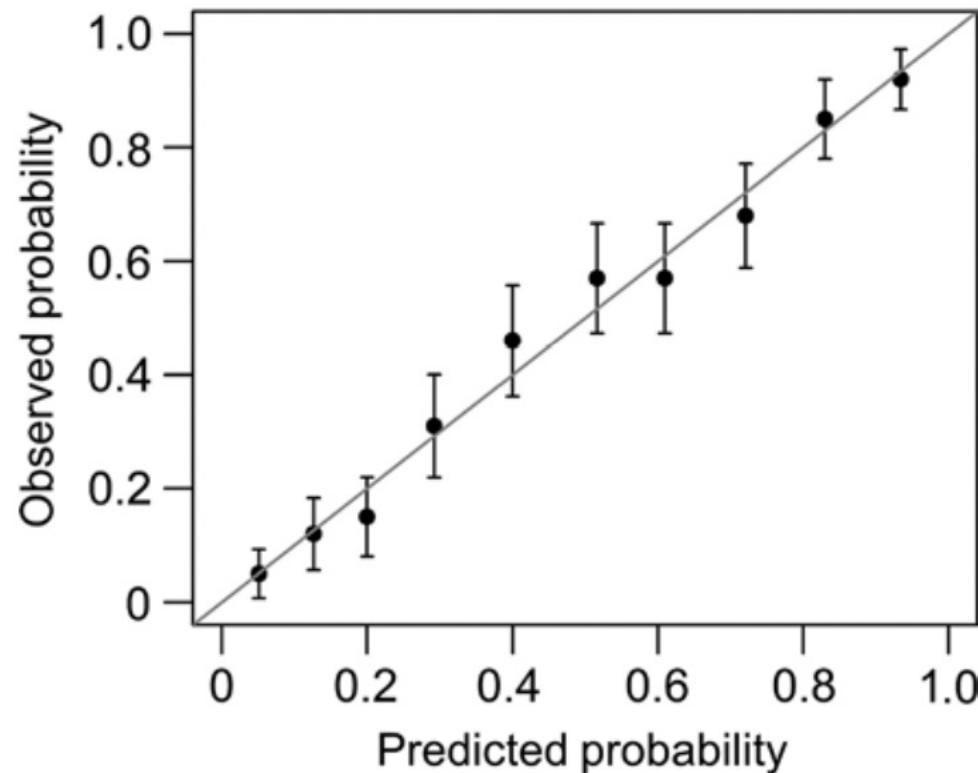
- Comparing segmented regions
 - Overlap measures
 - Dice coefficient
 - Intersection over union (Jaccard)
 - Surface distance
 - Hausdorff distance
- Evaluating clustering algorithms
 - Purity
 - Rand index
 - <http://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-clustering-1.html>

$$\text{Dice}(P, T) = \frac{|P_1 \wedge T_1|}{(|P_1| + |T_1|)/2}$$

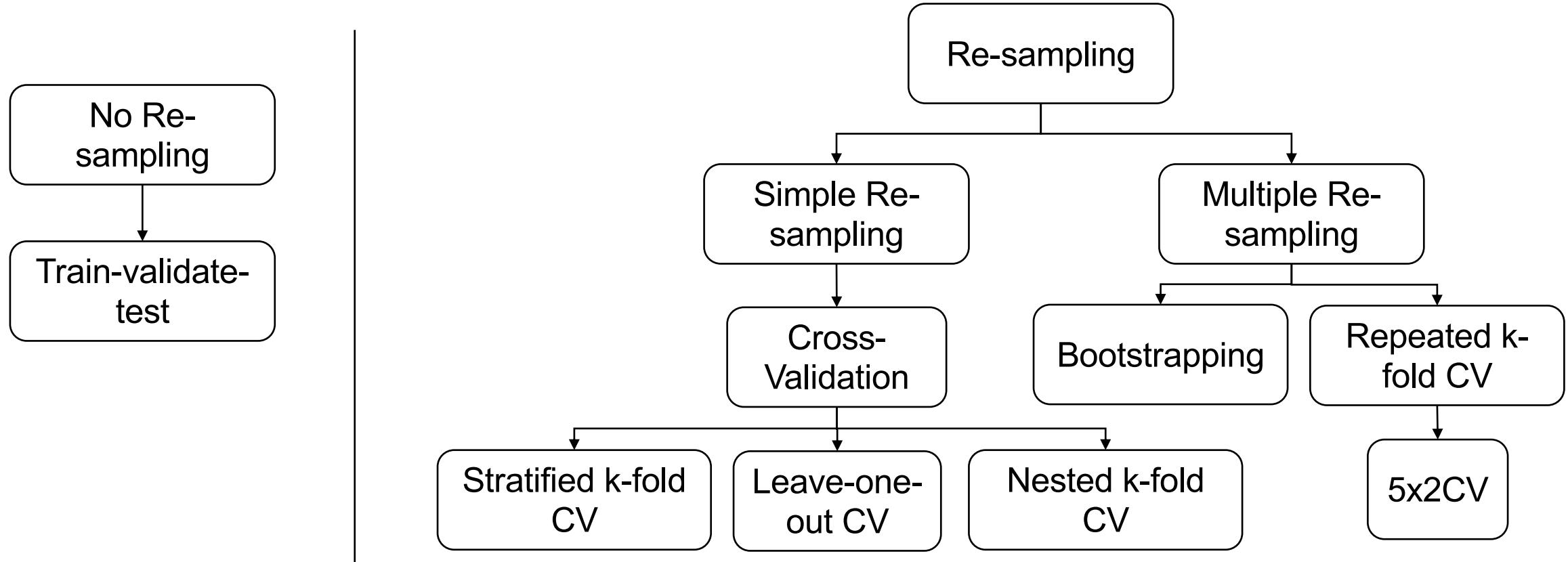


Calibration

- How well does the model's predicted probability approximates the actual event probability

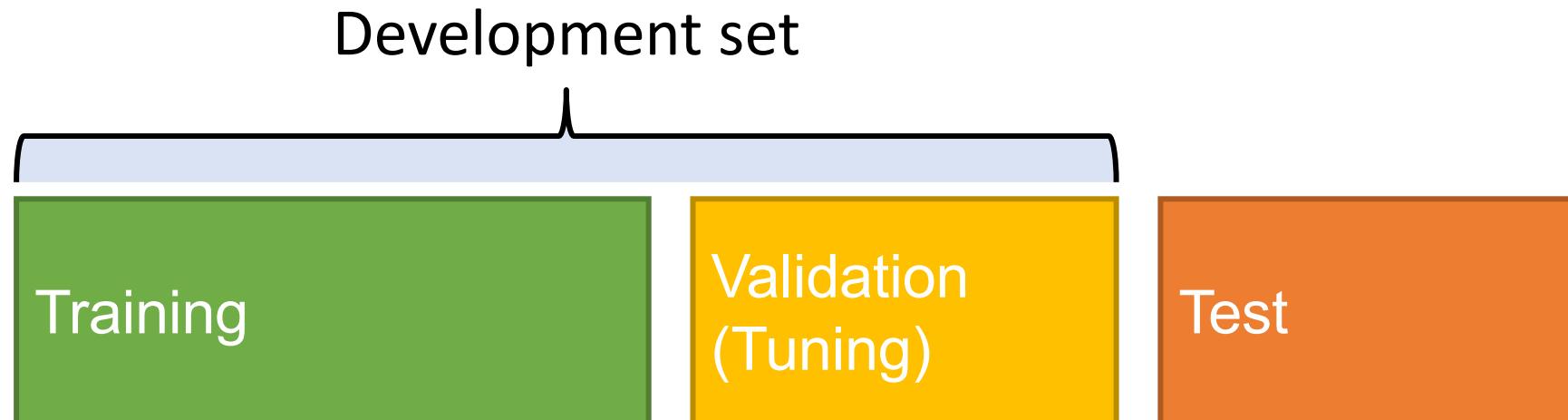


Evaluation Designs



Train-validation-test

- Split data into three groups
- Use training data to fit different models
- Use validation data to estimate generalization error (and tune model)
- Use test data to assess the performance

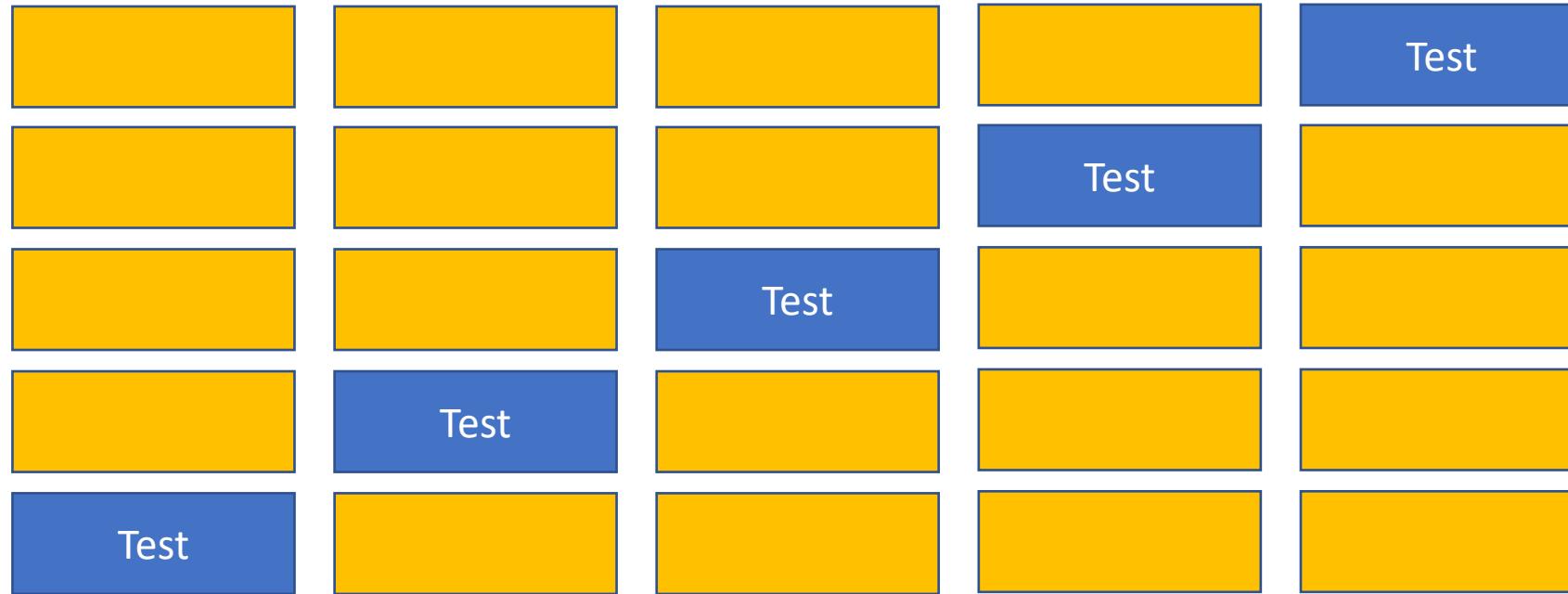


Notes on Train-Validate-Test

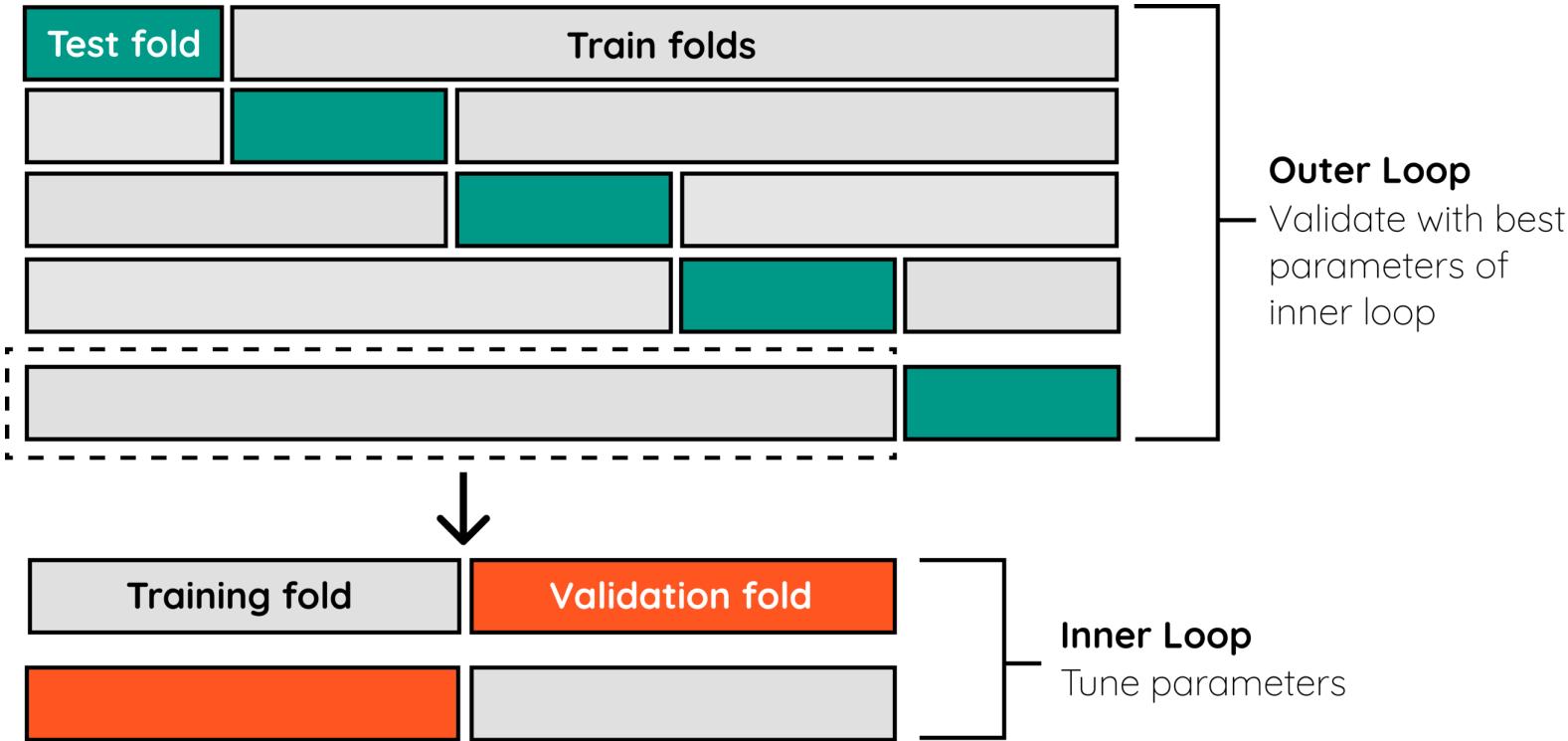
- Often need to balance between number of training and testing cases
 - Larger training set provides the learning algorithm with a more representative sample of the data
 - Larger test sets provide a more reliable estimate of accuracy (lower variance estimate)
- A single training set does not reveal how sensitive the accuracy of the learning algorithm is to a particular training sample → sampling bias

Cross validation

- Partition data into k-subsamples then iteratively assign one subsample as the test set, train the learning algorithm on the rest

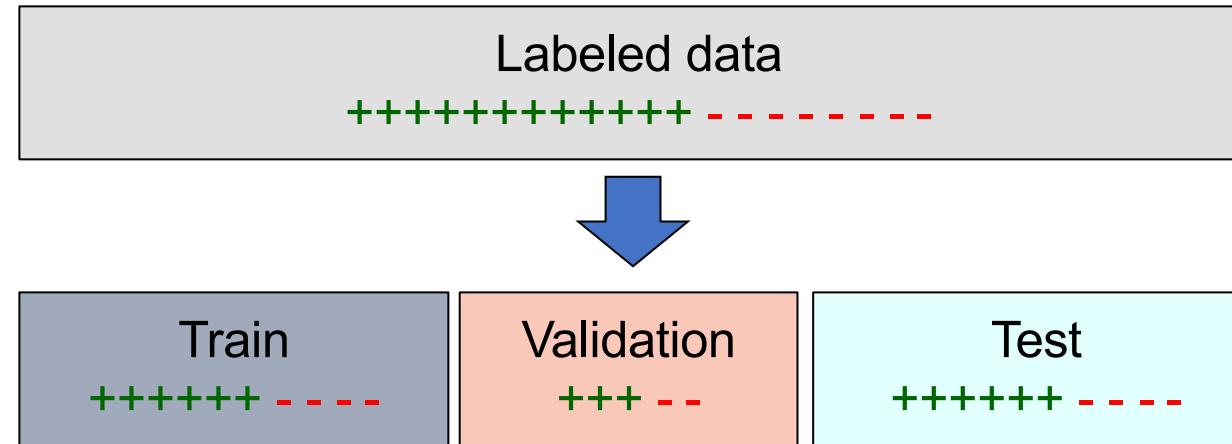


Nested Cross Validation



Source: <https://towardsdatascience.com/validating-your-machine-learning-model-25b4c8643fb7>

Stratified Cross-Validation

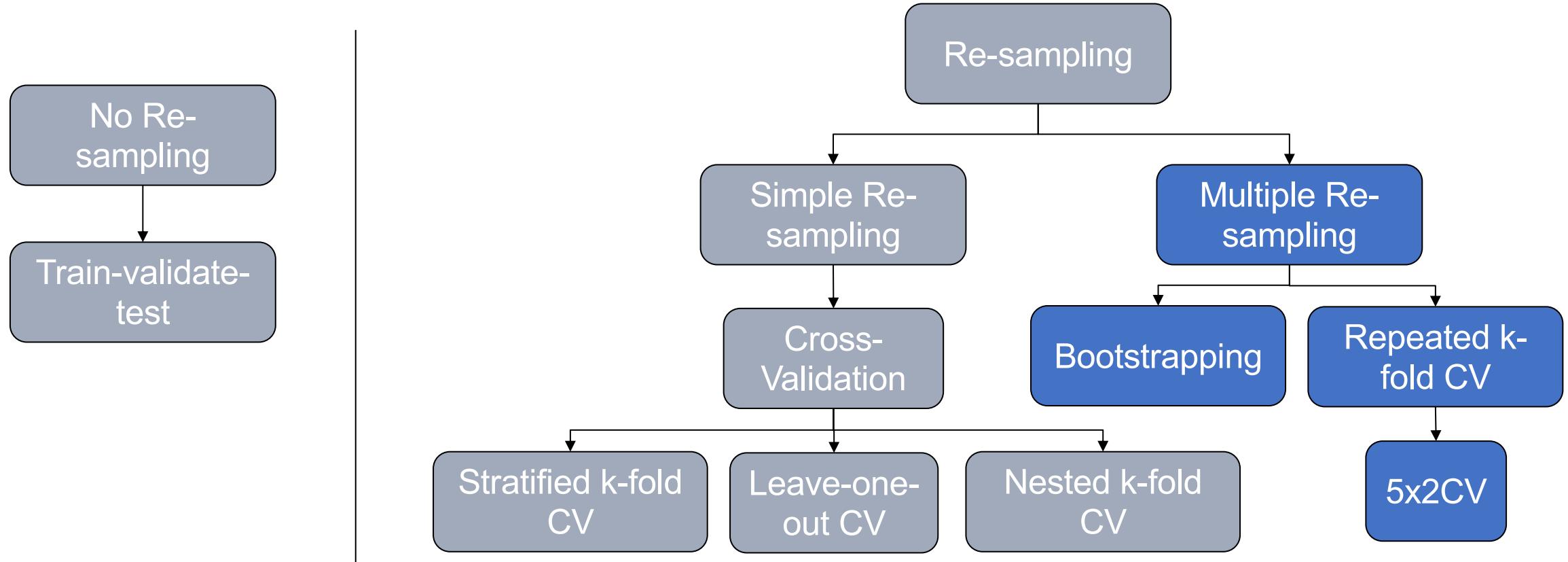


When randomly selecting training or validation sets, stratified cross-validation is used to ensure that class proportions are maintained in each selected set

Notes about cross validation

- The goal is to evaluate the learning algorithm, not an individual learned model
- 10-fold cross validation is commonly used, but smaller values of k can be used when learning takes a lot of time
- When working with time series data, training data **must** always be selected **before** the test data (otherwise training data will contain information from the future)

Resampling Methods



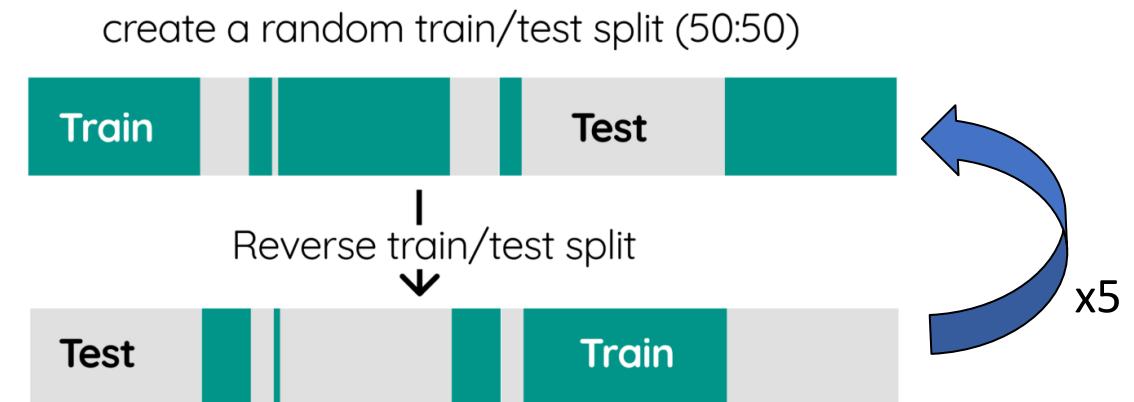
Adapted from: Japkowicz and Shah. <https://doi.org/10.1017/CBO9780511921803.004>

Bootstrapping

- Bootstrapping assumes that available samples are representative and creates a large number of new samples by drawing with replacement from the available samples
- Bootstrapping can be useful when the sample size is too small for cross validation or leave-one-out approaches to yield a good estimate

Repeated k-fold CV

- To obtain more stable estimates of an algorithm's performance, perform multiple runs of simple re-sampling schemes
 - Dietterich (1998) proposed the 5x2CV, in which a two-fold cross-validation is repeated 5 times



Source: <https://towardsdatascience.com/validating-your-machine-learning-model-25b4c8643fb7>

Comparing Models: Motivation

- Which model is better?

Algo	Acc	RMSE	TPR	FPR	Prec	Rec	F	AUC	Info S
NB	71.7	.4534	.44	.16	.53	.44	.48	.7	48.11
C4.5	75.5	.4324	.27	.04	.74	.27	.4	.59	34.28
3NN	72.4	.5101	.32	.1	.56	.32	.41	.63	43.37
Ripp	71	.4494	.37	.14	.52	.37	.43	.6	22.34
SVM	69.6	.5515	.33	.15	.48	.33	.39	.59	54.89
Bagg	67.8	.4518	.17	.1	.4	.17	.23	.63	11.30
Boost	70.3	.4329	.42	.18	.5	.42	.46	.7	34.48
RanF	69.23	.47	.33	.15	.48	.33	.39	.63	20.78

Source: N Japkowicz. https://www.icmla-conference.org/icmla11/PE_Tutorial.pdf

What is being tested?

Are the performance measurements attributable to the machine learning algorithm or are they observed by chance?

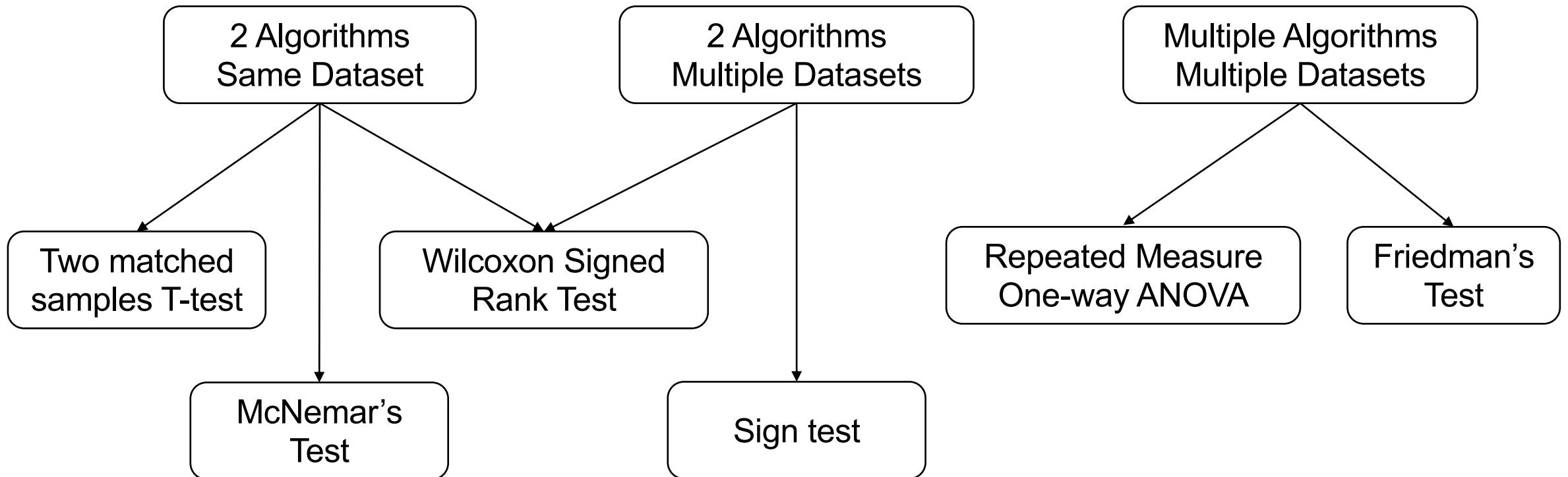
- Provide evidence that our selected evaluation measure is representative of the general behavior of our machine learning algorithm
- Parametric versus nonparametric test
 - Parametric tests make strong assumptions about the distribution of the underlying data (normal, Bernoulli, etc) but require smaller sample size
 - Nonparametric tests make weaker assumptions (skewed distribution, presence of outliers) but are also less powerful (less apt at rejecting the null hypothesis when it is false)

Considerations

- Can a parametric versus nonparametric test be used
- Select appropriate test based on the following scenarios
 - The comparison of 2 algorithms on the same dataset
 - The comparison of 2 algorithms on different datasets
 - The comparison of multiple algorithms on multiple datasets

Reference: Japkowicz and Shah. <https://doi.org/10.1017/CBO9780511921803.004>

Test Selection



The paired t-test

- Given two matched samples
 - the results of two classifiers applied to the same data set with **matching randomizations and partitions**
- Test whether the difference in **mean of a metric (e.g., accuracy)** between these two classifiers is significant
 - Test whether the two samples come from the same population.
 - Look at the distribution in **observed means and standard deviation**.
- We assume that the difference between means is zero (the null hypothesis) and see if we can reject this hypothesis.

Assumptions of the t-test

- The **normality** or pseudo-normality assumption: The t-test requires that the samples come from normally distributed population.
- The **randomness** of the samples: The sample should be representative of the underlying population. Therefore, the instances of the testing set should be randomly chosen from their underlying distribution.
- **Equal variance** of the populations: The two sample come from populations with equal variance.

McNemar's test

- Non-parametric counterpart of the t-test observed on the test set
 - # of instances misclassified by both classifiers C_{00}
 - # of instances misclassified by first classifier but correctly classified by the second classifier C_{01}
 - # of instances misclassified by the second classifier but correctly classified by the first classifier C_{10}
 - # of instances correctly classified by both classifiers C_{11}

$$X_{MC}^2 = \frac{(|C_{01} - C_{10}| - 1)^2}{C_{01} + C_{10}}$$

Two classifiers, multiple domains

- There are no clear parametric way to deal with the problem of comparing the performance of two classifiers on multiple domains:
 - The t-test is not a very good alternative because it is not clear that we have commensurability of the performance measures across multiple domains
 - The normality assumption is difficult to establish, e.g., the number of datasets tested > 30
- The t-test is susceptible to outliers, which is more likely when many different domains are considered
- Two non-parametric alternatives
 - The Sign Test
 - Wilcoxon's signed-Rank Test

The Sign test

- The sign test can be used either to compare
 - two classifiers on a single domain (using the results at each fold as a trial)
 - two classifiers on multiple domains
- We count
 - number of times that algorithm 1 outperforms algorithm 2, n_{a1}
 - number of times that algorithm 2 outperforms algorithm 1, n_{a2}
- The null hypothesis (stating that the two classifiers perform equally well) holds if the number of wins follows a binomial distribution.
- Practically speaking, a classifier should perform better on at least w_α datasets to be considered statistically significantly better at the α significance level, where w_α is the critical value for the sign test at the α significance level

Wilcoxon's signed-Rank Test

- Non-parametric
- For each domain, we calculate the difference in the performance of the two classifiers.
 - Rank the absolute values of these differences and graft the signs in front of the ranks.
 - Calculate the sum of positive and negative ranks, respectively (W_{S1} and W_{S2})
 - $T_{\text{Wilcox}} = \min(W_{S1}, W_{S2})$
- Compare to critical value V_α . If $V_\alpha \geq T_{\text{Wilcox}}$ we reject the null hypothesis that the performance of the two classifiers is the same, at the α confidence level.

Notes about statistical testing

- Keep in mind the assumptions set by the selected statistical test (e.g., normality, equal variance)
- Failure to achieve a statistically significant result cannot be interpreted as a true lack of difference – especially when the study is statistically underpowered
 - Statistical significance does not imply that the result is clinically meaningful
- When multiple statistical hypotheses are tested using the same dataset, the chance of observing a rare event increase, increasing the likelihood of incorrectly concluding that a real effect has been observed

Common Study Critiques

Lack of information

- Inconsistent reporting of datasets, models, data processing, and annotations → **recommendations for reporting**

Definition of the "gold standard"

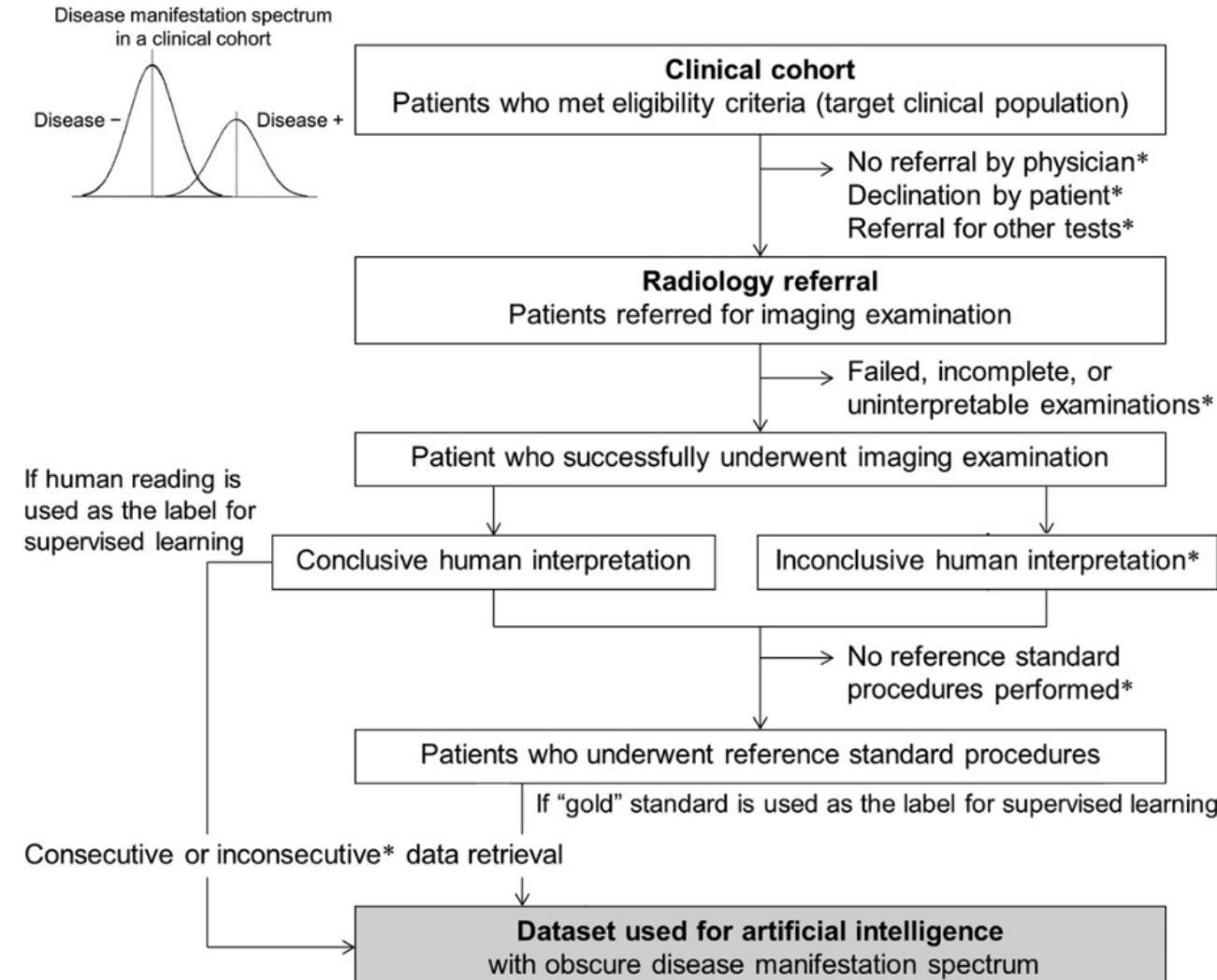
- Clarity about how “truth” was determined → **recommendations for providing baseline comparisons**

Appropriate sample size?

- Uncertainty as to how much training/testing examples are "enough"? → **recommendations for defining study population**

Reference: Gregory J et al, J Mag Res Imaging. <https://doi.org/10.1002/jmri.27035>

Understanding the cohort



Defining a reference standard

- Source of the standard
 - Completeness?
 - Consistent?
- Rationale for choosing the standard
- How was the standard defined?
 - Subject to measurement error?
 - Subjective judgment?

		Majority Decision of Retinal Specialist Grading before Adjudication				
		No	Mild	Moderate	Severe	Proliferative
Adjudicated Consensus	No	1469	4	5	0	0
	Mild	58	62	5	0	0
	Moderate	22	3	118	1	0
	Severe	0	0	13	36	1
	Proliferative	0	0	0	1	15

Confusion matrix for diabetic retinopathy between the grade determined by majority decision and adjudicated consensus.

Source: <https://doi.org/10.1016/j.ophtha.2018.01.034>

Guidelines for reporting data

- Data sources
 - Inclusion criteria
 - Selection of data subsets
- Data pre-processing steps
 - Definition of data elements
 - De-identification methods
 - How missing data were handled
- Reference standard
 - How was the standard defined?
 - The rationale for choosing the standard
 - Source of ground truth annotations
 - If generated by human annotators, were inter and intra-rater variability mitigated?

Guidelines for reporting models

- Model
 - Description of model
 - Inputs and outputs
 - Hyperparameters
- Implementation
 - Software libraries and frameworks
 - Computational requirements
- Training
 - How was the dataset split for training?
 - Was transfer learning used?
 - Method of selecting final model

Guidelines for reporting evaluation

- Evaluation approach
 - Metrics used to measure model performance
 - Statistical tests used
 - Whether sensitivity analysis was performed to evaluate model robustness
 - Techniques for explaining model results, if applicable
 - Validation of model using external data
- Results
 - Whether performance metrics were compared with previously reported models, if possible
 - Whether confidence intervals are provided
 - Failure analysis
- Limitations
 - Possible sources of bias
 - Implications for practice

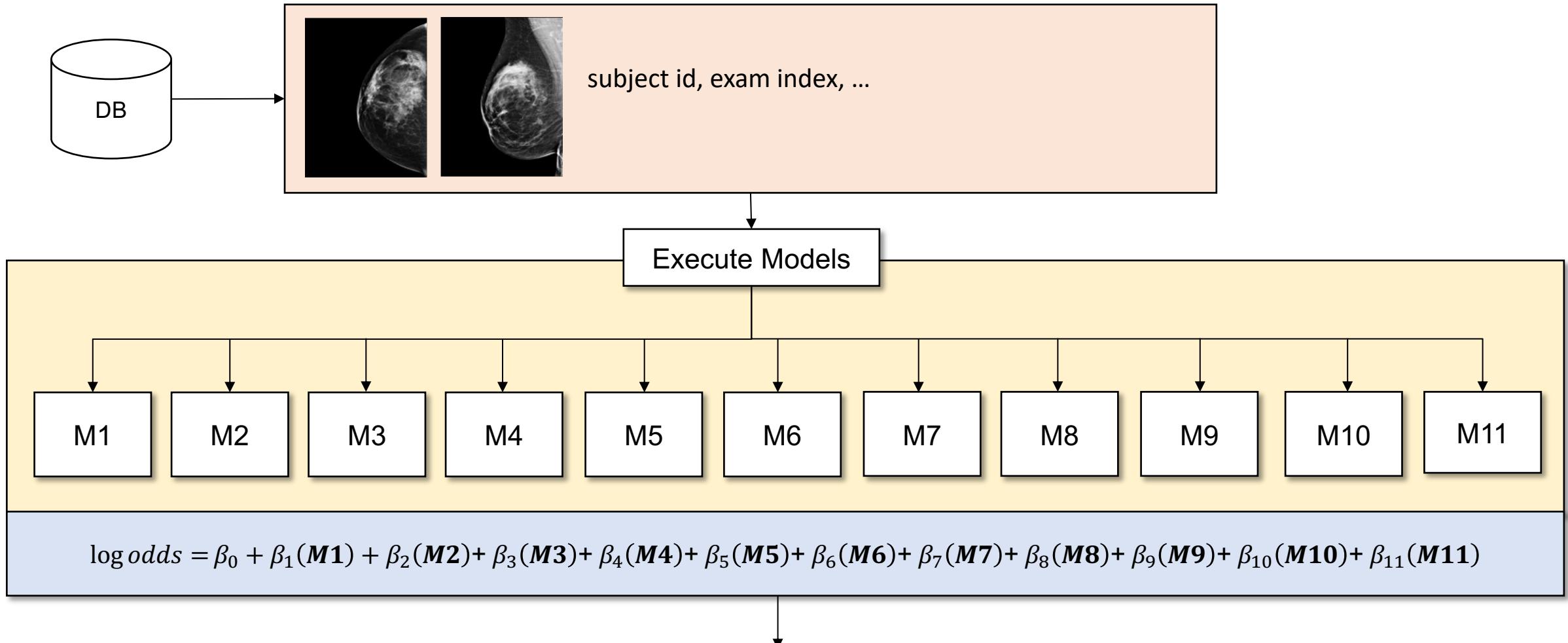
Helpful reporting standards

- Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD)
 - <https://www.equator-network.org/reporting-guidelines/tripod-statement/>
 - TRIPOD-AI under development: <https://osf.io/zyacb/>
- Prediction model Risk Of Bias ASsessment Tool (PROBAST)
 - <https://www.acpjournals.org/doi/10.7326/M18-1376>
 - PROBAST-AI under development
- Reporting of clinical trials involving AI
 - CONSORT/CONSORT-AI: <https://www.nature.com/articles/s41591-020-1034-x>
 - SPIRIT (protocol reporting): <https://www.nature.com/articles/s41591-020-1037-7>

Example: External validation of breast AI

- New breast AI/ML algorithms are becoming commercially available and are being reported to have performance on par with breast radiologists
 - Lack of transparency about the cohorts in which the algorithm was trained and evaluated
→ will the performance generalize?
 - Guidelines on how institutions should review and adopt AI/ML algorithms are sparse
→ what factors should institutions consider before buying?
 - Check if the AI/ML algorithm works as intended
→ when does an algorithm deviate from stated performance?
- Objective: Perform an external validation of the "DREAM Challenge Competition Ensemble Method" model on a retrospective cohort from our institution
 - Primary measures: AUC, sensitivity, specificity
 - Identify conditions when model performance differs from previously reported results

Competition Ensemble Method model

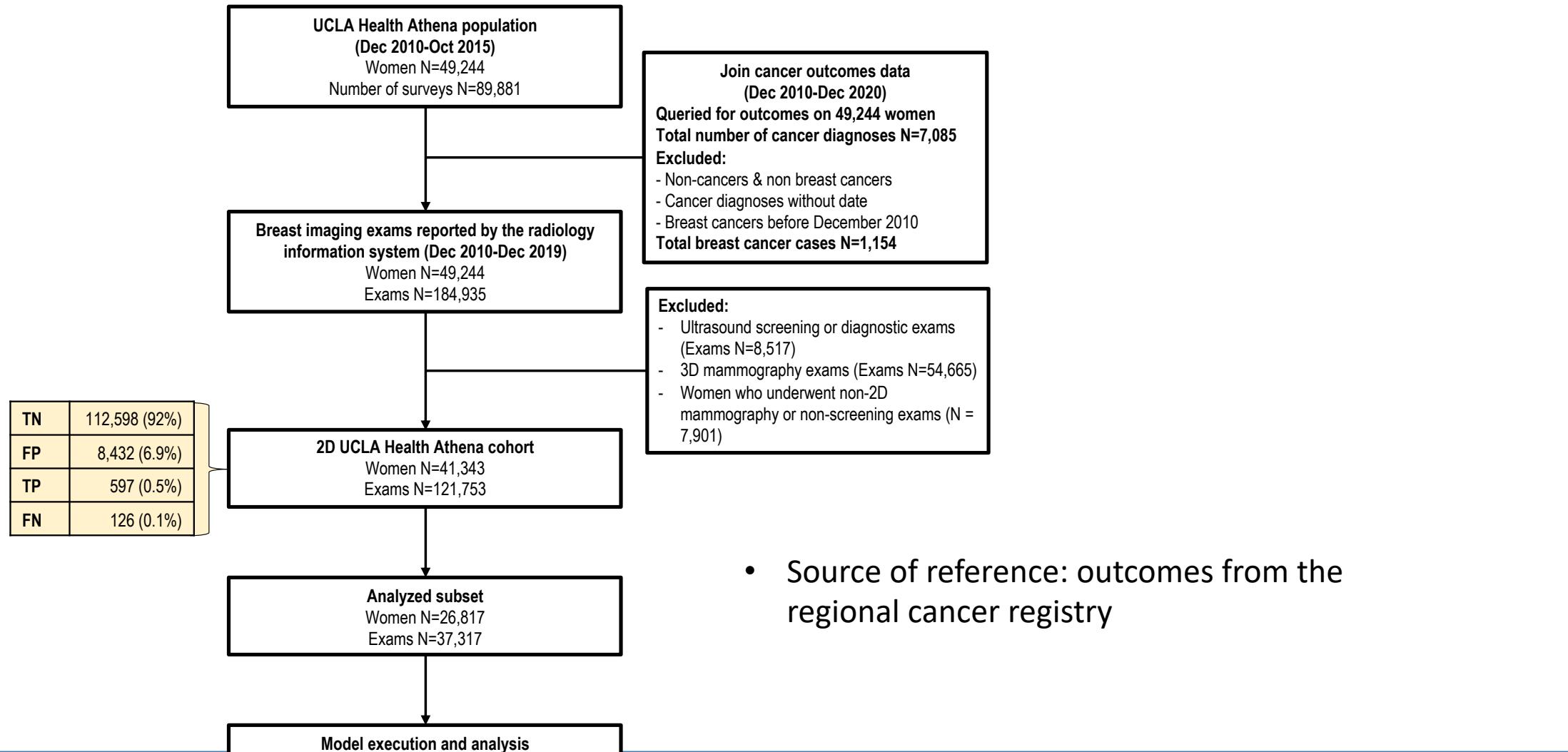


Overall score: 0.871 → Cancer

SPIE. MEDICAL IMAGING

Reference
UCLA Health
David Geffen School of Medicine

External validation: Cohort definition



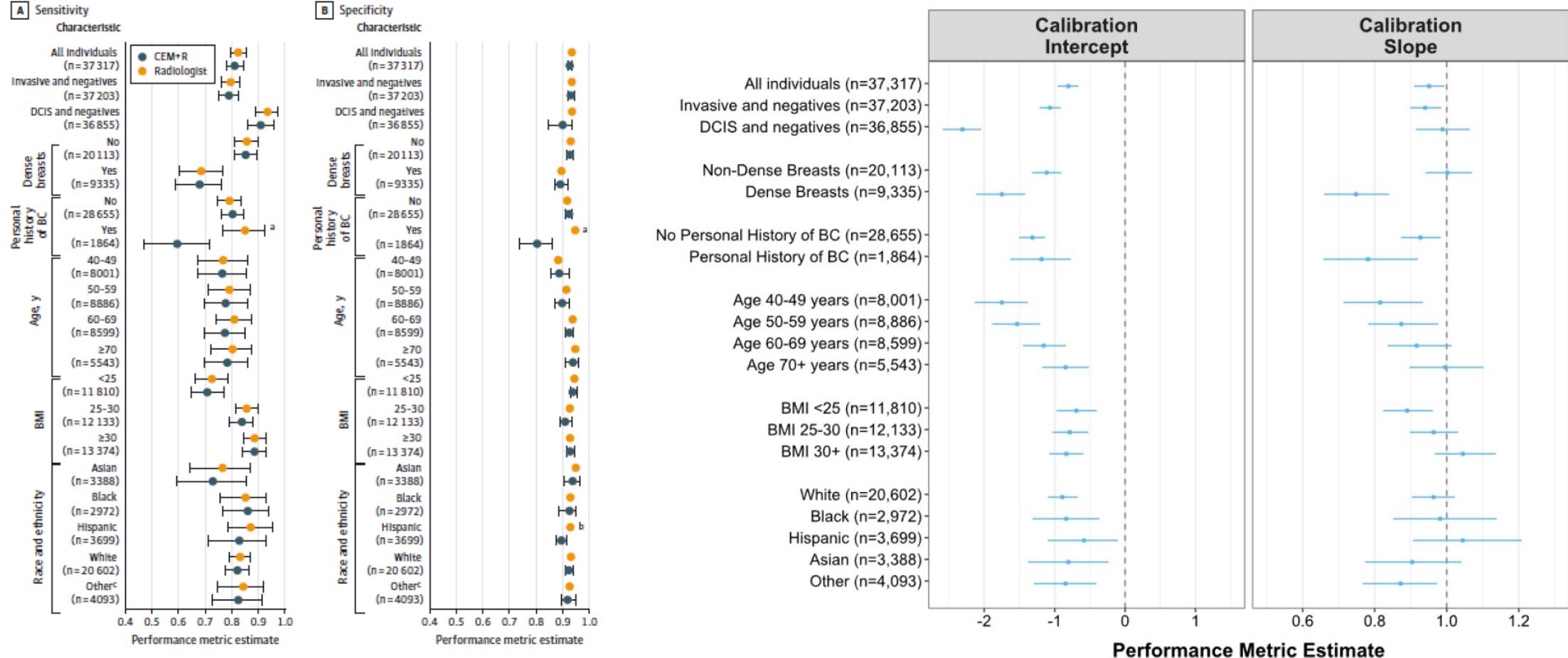
Source: <http://jamanetwork.com/article.aspx?doi=10.1001/jamanetworkopen.2022.42343>

SPIE. MEDICAL IMAGING

UCLA Health

David Geffen School of Medicine

Results: Subgroup analysis



High-performance medicine: the convergence of human and artificial intelligence

Eric J. Topol 

The use of artificial intelligence, and the deep-learning subtype in particular, has been enabled by the use of labeled big data, along with machine learning, at the core of AI applications in medicine. This is beginning to have important clinical applications.

Validation of the performance of an algorithm in terms of its accuracy is not equivalent to demonstrating clinical efficacy. This is what Pearse Keane and I have referred to as the ‘AI chasm’—that is, an algorithm with an AUC of 0.99 is not worth very much if it is not proven to improve clinical outcomes.

Resources

- Latest course materials
 - <https://uclawillhsu.github.io/spie2023mi/>
- Evaluating learning algorithms (Japkowicz and Shah)
 - <https://doi.org/10.1017/CBO9780511921803>
- Chapter 10: Evaluation (Medical Imaging Informatics, Bui and Taira, eds. 2010)
 - https://link.springer.com/chapter/10.1007/978-1-4419-0385-3_10

Acknowledgements and references



William Hsu, PhD

Associate Professor of Radiological Sciences, Bioinformatics, and Bioengineering

wlsu@mednet.ucla.edu
<http://hsu-lab.com>

Research team (Hsu Lab @ Medical & Imaging Informatics)



The Hsu Lab gratefully acknowledges support from the National Science Foundation (#1722516), National Institutes of Health (R01 EB031993; R01 CA271034; R01 HL127153; R01 CA22079; R01 EB029346; R01 CA210360; R37 CA240403), UCLA SPORE in Prostate Cancer, Jonsson Comprehensive Cancer Center, and the Department of Radiological Sciences.

Brief read

- Liu Y, Chen PH, Krause J, Peng L. How to read articles that use machine learning: users' guides to the medical literature. *Jama*. 2019 Nov 12;322(18):1806-16.

More in-depth

- Steyerberg, EW. Evaluation of Performance. In: Clinical Prediction Models. Statistics for Biology and Health. Springer, Cham; 2019.

Textbook reference

- Japkowicz N, Shah M. Evaluating learning algorithms: a classification perspective. Cambridge University Press; 2011.