

Lecture 6: Assessment of machine (supervised) learning systems



Universitat
de les Illes Balears

Departament
de Ciències Matemàtiques
i Informàtica

11752 Aprendizaje Automático
11752 Machine Learning
Máster Universitario
en Sistemas Inteligentes

Alberto ORTIZ RODRÍGUEZ

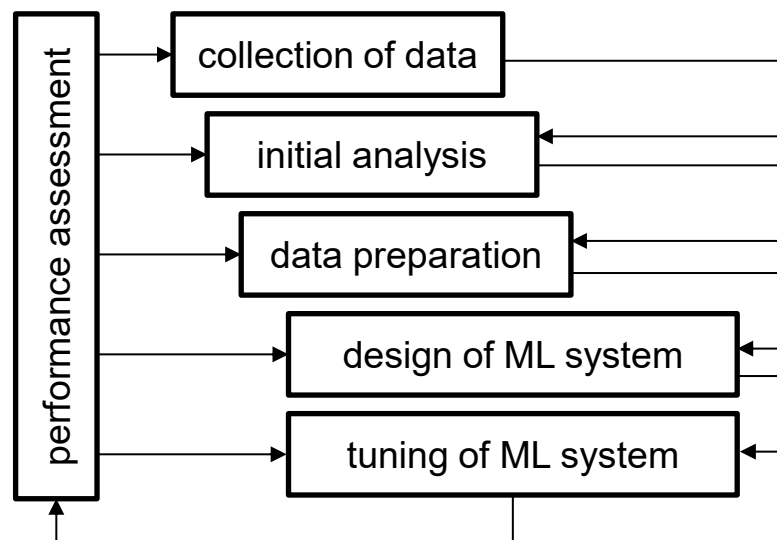
- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

- The assessment of ML systems involves several aspects of its performance and may take the designer back to any of the developing stages
- Regarding the model itself, we may evaluate

underfitting
vs
overfitting

- its **adequacy**
 - is the model adequate? is it biased?
which is the error rate?
- its **generalization**
 - how well behaves the model with unseen data

- ML systems evaluation tools are useful not only for assessing the performance of the system, but also
 - for **debugging** purposes, to figure out what is not working and make the necessary adjustments
 - to **compare among different ML models**
- The basic assessment protocol splits the dataset into:
 - **Training set**. Used to build the classifier/regressor.
 - **Test set**. To check the behaviour of the classifier with unseen examples.



- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

Bias-variance tradeoff

- The source of prediction errors can be shown to be three-fold:

- **Bias**

- It is due to wrong assumptions, either by the designer or by the model
- A high-bias model is most likely to **underfit** the training data

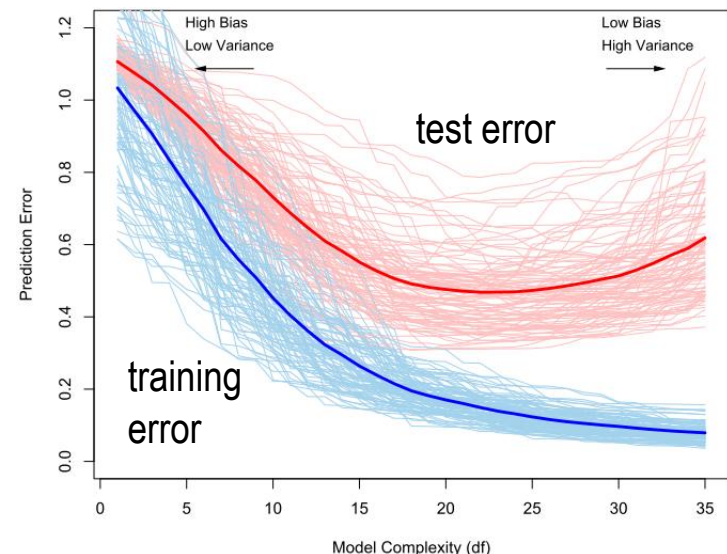
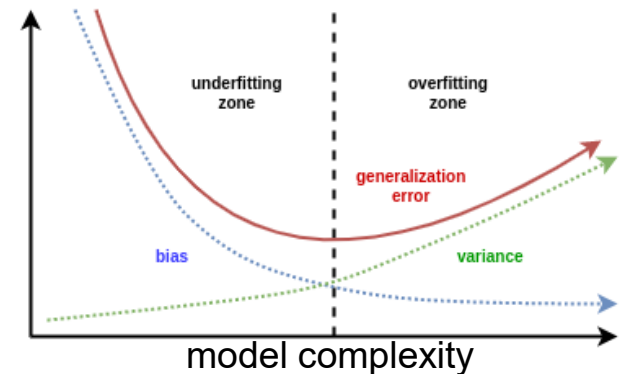
- **Variance**

- It is due to the model's **excessive sensitivity** to small variations in the training data
- A model with many degrees of freedom, e.g. a high-degree polynomial, is likely to have high variance and thus to **overfit** the training data

- **Irreducible error**

- It is due to the **noisiness** of the data itself
- the only way to reduce this error is to clean up the data, i.e. fix a broken sensor

- **Increasing a model's complexity** typically increases its variance and reduces its bias
- **Reducing a model's complexity** usually increases its bias and reduces its variance
- This is why it is called the **bias-variance tradeoff** (also BV dilemma)



- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

Confusion matrix

- Many **performance metrics** can be calculated from the **confusion matrix**, e.g. for **two** classes, having defined first which is the positive class

		predicted class		
		positive	negative	
true class	positive	TP	FN	$P_o = TP + FN$
	negative	FP	TN	$N_e = FP + TN$
$\widehat{P_o} = TP + FP \quad \widehat{N_e} = FN + TN$				

accuracy (A)	$\frac{TP+TN}{TP+TN+FP+FN}$	error rate (E)	$1 - A = \frac{FP+FN}{TP+TN+FP+FN}$
false positive rate (FPR)	$\frac{FP}{N_e} = \frac{FP}{FP+TN}$	true positive rate (TPR)	$\frac{TP}{P_o} = \frac{TP}{TP+FN}$
precision (P)	$\frac{TP}{\widehat{P_o}} = \frac{TP}{TP+FP}$	recall (R)	$\frac{TP}{P_o} = \frac{TP}{TP+FN} = TPR$
specificity	$\frac{TN}{N_e} = 1 - FPR$	F_1 -score	$\frac{2}{\frac{1}{P} + \frac{1}{R}} = \frac{2PR}{P+R}$

FPR is also known as **false alarm rate**

TPR is also known as **sensitivity**

- **F₁ -score** combines in a single metric P and R, by means of their **harmonic mean**

$$\frac{2}{1/P + 1/R} = \frac{2PR}{P + R}$$

- the regular mean treats all values equally, the harmonic mean penalizes more low values
- a high F₁ score (which is better) results only if both P and R are high
- also known as Sorensen-Dice coefficient or **Dice Similarity Coefficient** (DSC)

- F₁ is a particular case of the **F - score** or **F_β - score**:

$$F_{\beta} = (1 + \beta^2) \frac{P R}{\beta^2 P + R} = \frac{(1 + \beta^2) TP}{(1 + \beta^2) TP + \beta^2 FN + FP}$$

$$F_1 = \frac{2TP}{2TP + FN + FP}, \quad F_2 = \frac{5TP}{5TP + 4FN + FP}, \quad F_{0.5} = \frac{1.25 TP}{1.25 TP + 0.25 FN + FP}$$

- β = 1, F₁ -score, which weighs equally R and P: same emphasis on FN and FP
- β = 2, **F₂ -score**, which weighs R lower than P: effect of FP is less noticeable
- β = 0.5, **F_{0.5} -score**, which weighs R higher than P: effect of FN is less noticeable

$$P = \frac{TP}{\widehat{P_o}} = \frac{TP}{TP + FP}$$

$$R = \frac{TP}{P_o} = \frac{TP}{TP + FN}$$

Confusion matrix

- The confusion matrix can be generalized for **M-class problems**

Confusion matrix

	en	es	fr	ge	it	ru
Actual	1819	3	51	69	57	1
	0	1742	0	0	0	258
fr	57	0	1810	36	95	2
ge	83	0	54	1831	32	0
it	58	5	89	20	1826	2
ru	0	331	0	0	0	1669
	en	es	fr	ge	it	ru
	Predicted					

TP, TN, FP, FN are calculated according to **one versus all (OvA) classification**

predicted class

	$C_0 \dots C_{k-1}$	C_k	$C_{k+1} \dots C_m$
true class	C_0	C_k	C_m
C_0	TN	FP	TN
C_k	FN	TP	FN
C_m	TN	FP	TN

- Global metrics can be calculated following two approaches:
 - macro-averages**: average scores for each class
 - each class is weighed equally
 - micro-averages**: from individual TP, TN, FP, FN
 - each prediction is weighed equally

$$P_k = \frac{TP_k}{TP_k + FP_k}, \text{ etc.}$$

$$\text{e.g. } P_{\text{macro}} = \frac{P_1 + \dots + P_m}{m}$$

$$\text{e.g. } P_{\text{micro}} = \frac{TP_1 + \dots + TP_m}{TP_1 + \dots + TP_m + FP_1 + \dots + FP_m}$$

- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

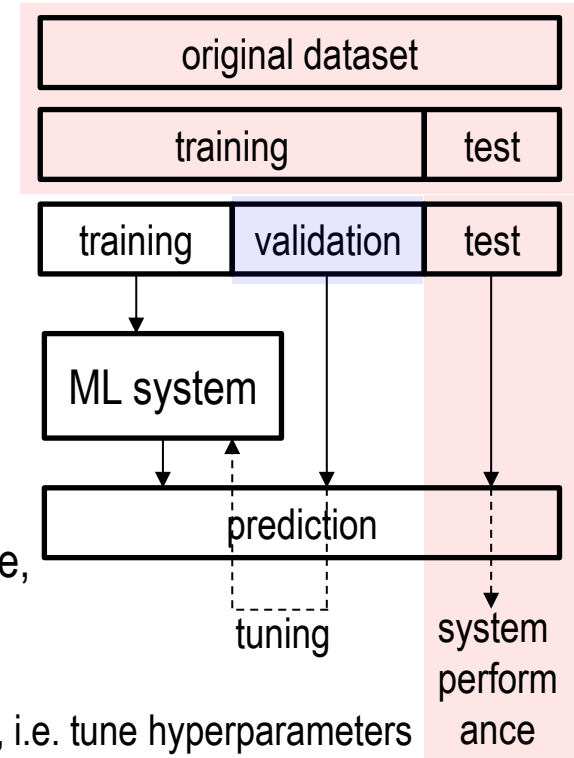
Cross-validation techniques

- From a **methodological** point of view, we need a way to obtain robust and bias-free performance measurements
- **Cross-validation** techniques can provide performance estimation values with low bias:
 - holdout cross validation
 - n-fold cross validation
 - stratified n-fold cross validation
 - leave-one-out cross validation (LOOCV)
 - nested cross validation

Cross-validation techniques

- **Holdout cross validation**

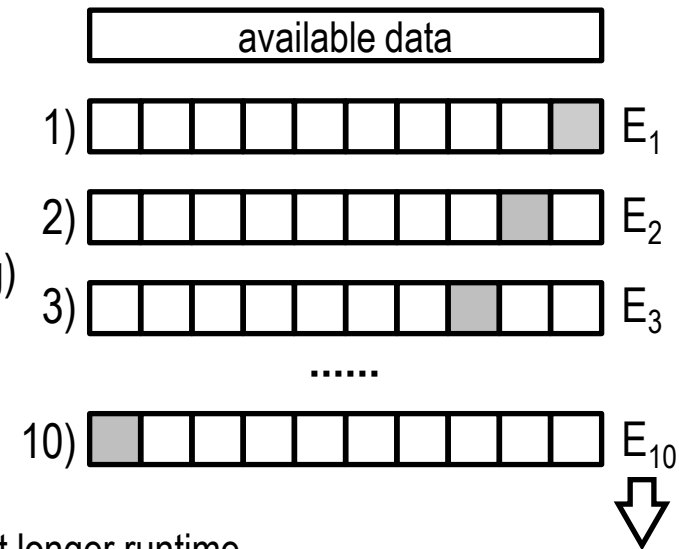
- Simplest kind of cross validation
- The data set is initially split into two sets: the **training set** and the **test set**
 1. ML system is built using the training set only
 2. ML system is asked to predict the output values for the data in the “unseen” test set
- The accumulated errors give the **test set error**, used to evaluate the model
 - Maybe high variance in the results if the test is repeated
- To tune the system appropriately and reduce this variance, we can split the training set in a **training subset** and a **validation subset**
 - The validation subset can be employed for model selection, i.e. tune hyperparameters
 - train the system
 - repeatedly evaluate it using the validation subset using different settings
- Even in this way, the evaluation may depend heavily on **which data points end up in the training set and which end up in the test set**, and thus the evaluation may be significantly different depending on how the splitting is made



Cross-validation techniques

- **n-fold cross validation**

- The available data is randomly split into n folds without replacement:
 - $n - 1$ folds are used for training
 - the remaining fold is used for testing
- The procedure is repeated n times, choosing a different fold for testing each iteration (no overlapping)
- A global performance measure is obtained by averaging the individual measurements
 - lower variance estimate than the holdout method
- In most cases, $n = 5$ or 10
 - n large \Rightarrow more data for training \Rightarrow lower bias in E but longer runtime
- A **high variance** indicates a situation of overfitting



$$E = \frac{1}{n} \sum_{i=1}^n E_i$$

$$V = \text{Var}(E_i)$$

- **Stratified n-fold cross validation**

- Class proportions are preserved in each fold
- Better performance estimates as for bias and variance

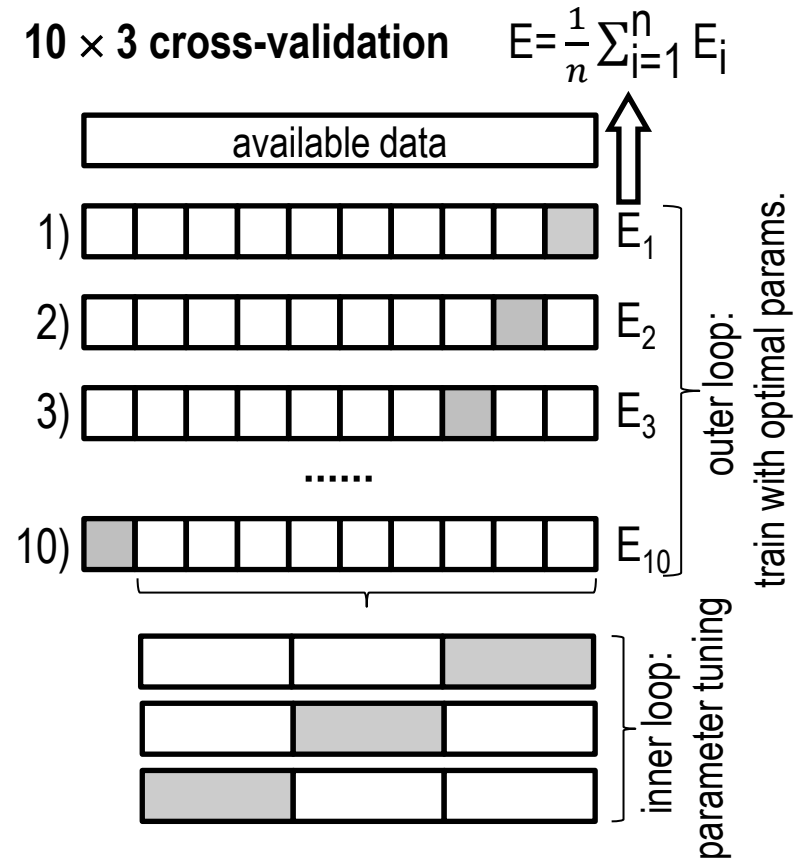
- **Leave-one-out cross validation (LOOCV)**

- $n = N$, the size of the dataset; hence, at each iteration, the test set comprises one single sample
- recommended for small datasets

Cross-validation techniques

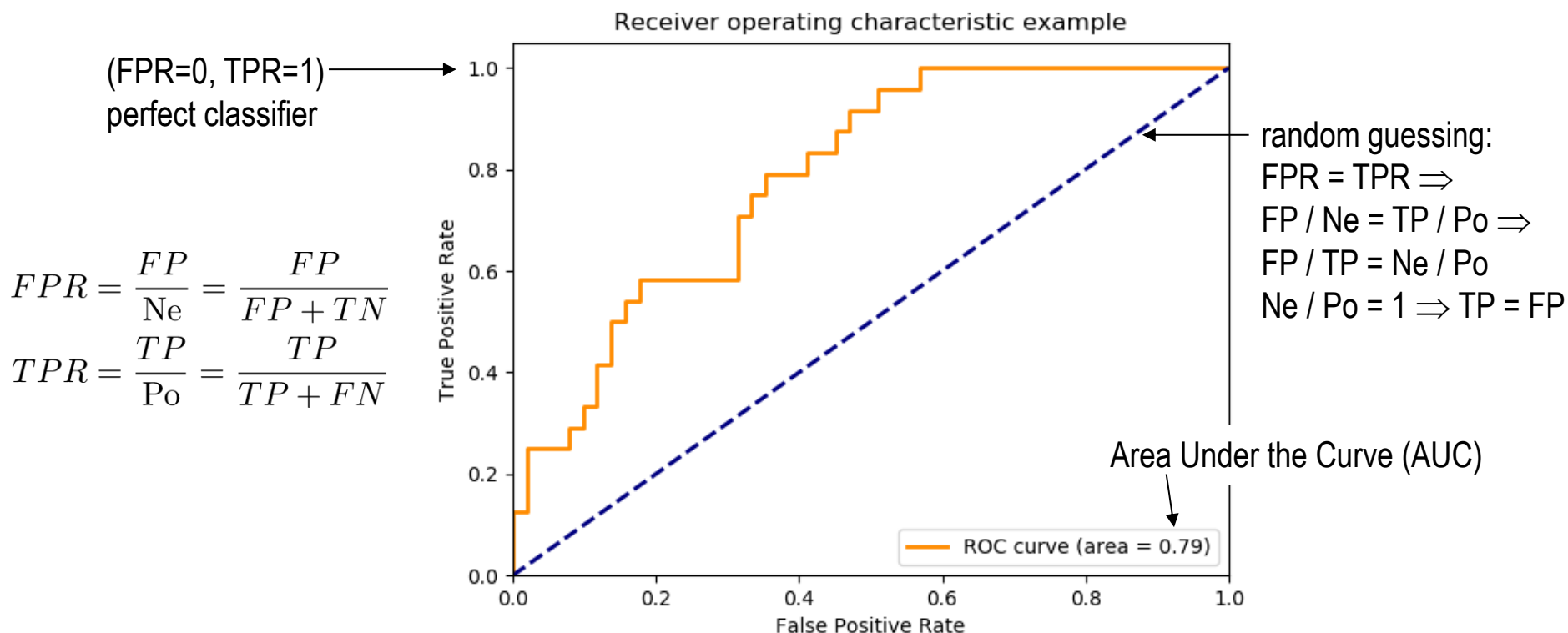
- **Nested cross validation**

- Outer n-fold: performance estimation
the available data is randomly split into **n folds** without replacement:
 - $n - 1$ folds are used for training
 - the remaining fold is used for test
- The procedure is repeated n times, choosing a different fold for testing each iteration
- Inner m-fold: model selection (hyperp. tuning)
the set of training folds is split into **m folds** leaving one for validation
- A global performance measure is obtained by averaging the individual measurements
- This measure gives a good estimate of what to expect from unseen data



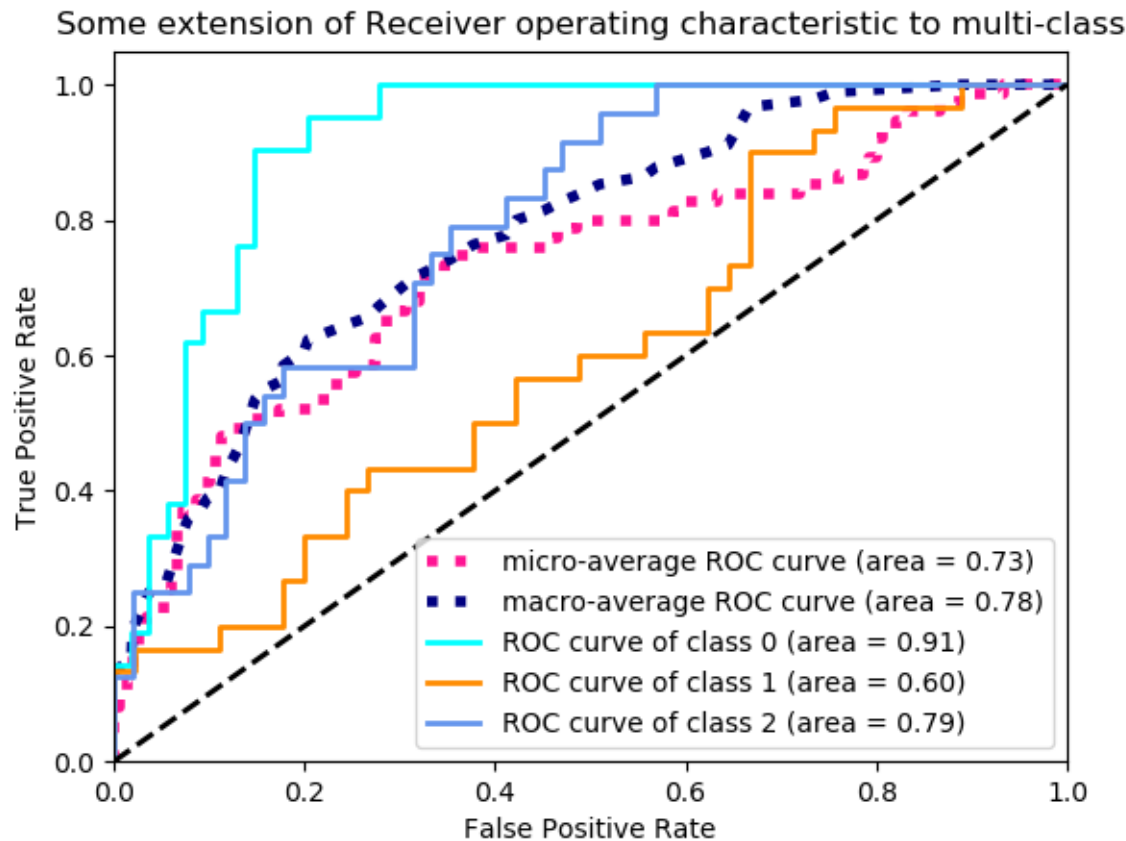
- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

- **Receiver Operating Characteristic** curve (concept from early Radar days)
 - Set of (FPR, TPR) points obtained by varying the algorithm's parameters
 - Every (FPR, TPR) point is 1 classifier configuration, choose the closest to the (0,1) point

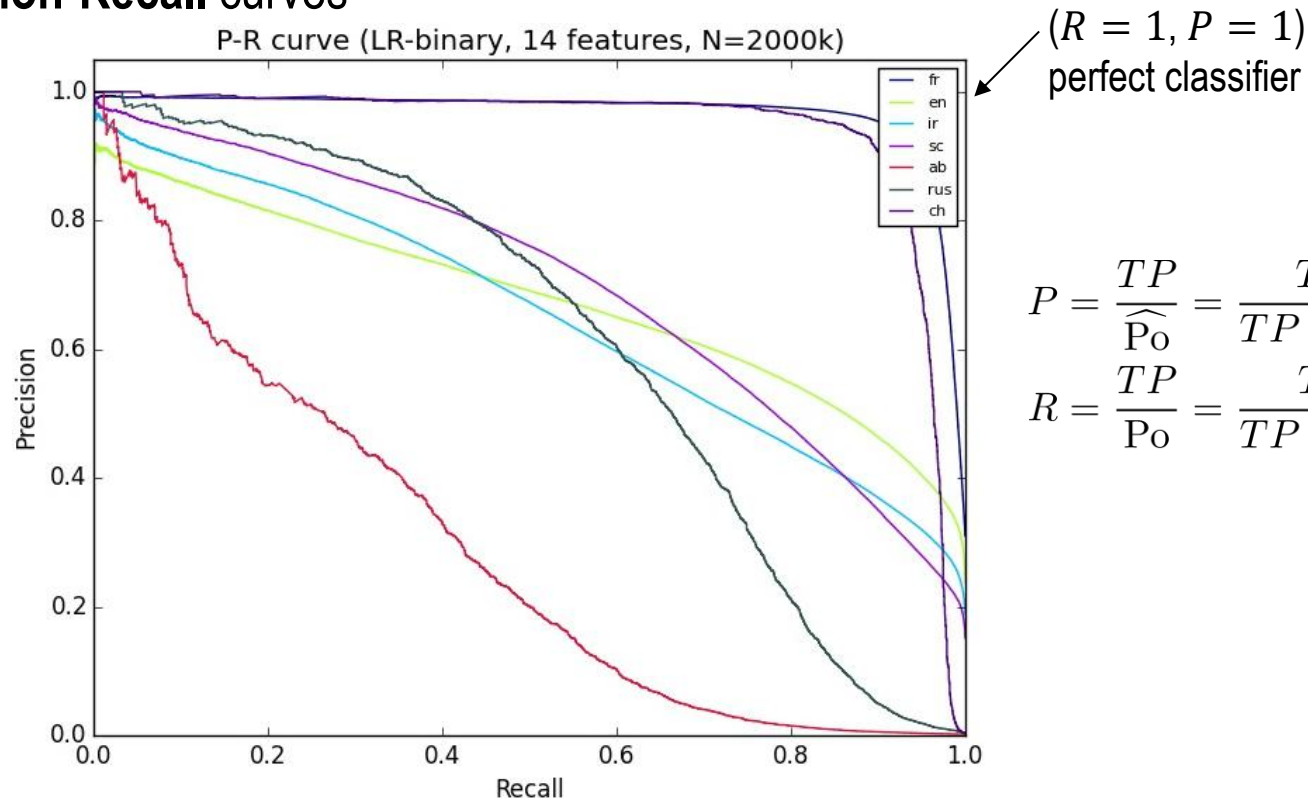


- Can also be shown as sensitivity (TPR) vs 1 – specificity (FPR)
 - **Area Under the Curve (AUC)**
 - Global measure of classifier performance, the higher the better

- **Receiver Operating Characteristic** curve (concept from early Radar days)
 - can also be calculated for multi-class problems



- Other usual curves for classifier performance characterization are the **Precision-Recall** curves



- The AUC is also of application in this case, also maximum R for P = 1, etc.
- Note: TN are not accounted for in this curve, it is the kind of problem where negatives are not as relevant as positives, e.g. inspection systems

- Introduction
- Bias-variance tradeoff
- Confusion matrix and performance metrics
- Assessment methodologies: Cross-validation techniques
- ROC curves
- Debugging and model tuning hints

Bias and variance problems

- By plotting the model training and validation accuracies as functions of the training set size (**learning curves**), we can easily detect whether the model suffers from high variance or high bias

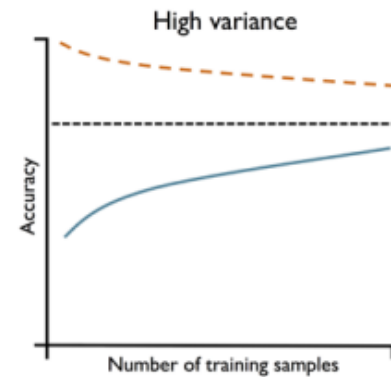
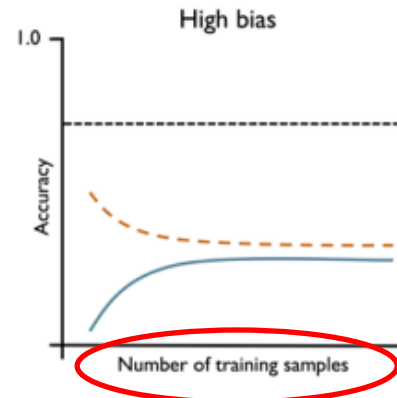
– high bias

- symptom: low training and validation accuracy and hence the model **underfits** the data
- solution: increase model complexity either with more parameters or more features

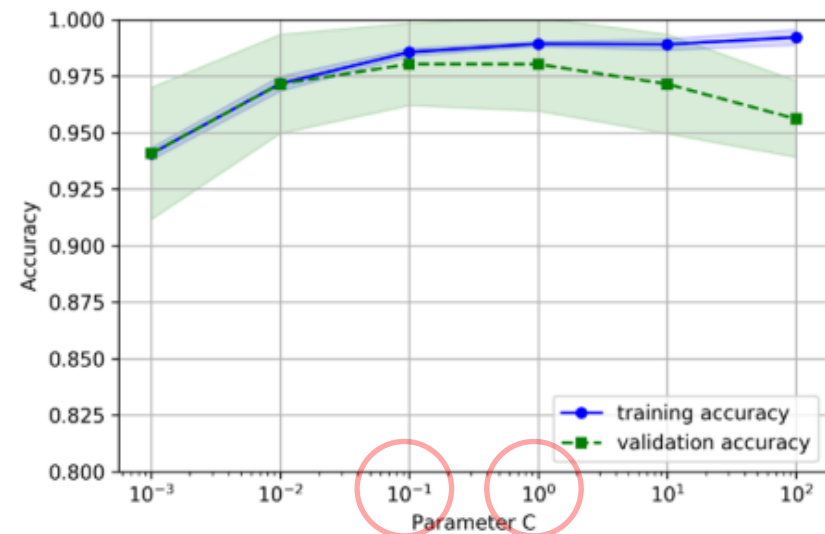
– high variance

- symptom: large gap between training and validation accuracy, the model captures too well the training data (**overfitting**)
- solution: collect more training data (care with noisy data, performance will not improve), reduce the model complexity (e.g. regularization) or remove some features (via feature selection or extraction)

progressively increase training set size, train and test



- In machine learning, we have two types of parameters:
 - those that are learned from training data, e.g. separating hyperplane
 - those that are optimized separately, e.g. decision tree depth or C in soft SVM
- The latter are known as **hyperparameters**
 - They can be set one by one, using **performance** (train and validation) **learning curves**



– or via **grid search**

- **brute-force** exhaustive search paradigm where we specify a list of values for different hyperparameters and the computer evaluates the model performance for each combination, e.g. when optimizing SVM

- kernel type: linear, (in)homogeneous polynomial, rbf
- kernel parameters: polynomial degree q, coefficient γ
- soft SVM hyperparameter C

$$K_{ln}(x, z) = x^T z$$

$$K_{hp}(x, z) = (x^T z)^q, \quad q > 0$$

$$K_{ip}(x, z) = (\gamma x^T z + 1)^q, \quad q > 0$$

$$K_{rbf}(x, z) = e^{-\gamma \|x - z\|^2}$$

$$\min J(w, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$

Lecture 6: Assessment of machine (supervised) learning systems



Universitat
de les Illes Balears

Departament
de Ciències Matemàtiques
i Informàtica

11752 Aprendizaje Automático
11752 Machine Learning
Máster Universitario
en Sistemas Inteligentes

Alberto ORTIZ RODRÍGUEZ