

Master Program in *Data Science and Business Informatics*

Statistics for Data Science

Lesson 30 - Classifier performances in R

Salvatore Ruggieri

Department of Computer Science

University of Pisa, Italy

salvatore.ruggieri@unipi.it

Tests and confidence intervals for classifier performance

The Caret package

- 1 Define sets of model parameter values to evaluate
- 2 **for** *each parameter set* **do**
- 3 | **for** *each resampling iteration* **do**
- 4 | | Hold-out specific samples
- 5 | | [Optional] Pre-process the data
- 6 | | Fit the model on the remainder
- 7 | | Predict the hold-out samples
- 8 | **end**
- 9 | Calculate the average performance across hold-out predictions
- 10 **end**
- 11 Determine the optimal parameter set
- 12 Fit the final model to all the training data using the optimal parameter set

For resampling methods, see Lesson 28

See R script

Binary classifier performance metrics

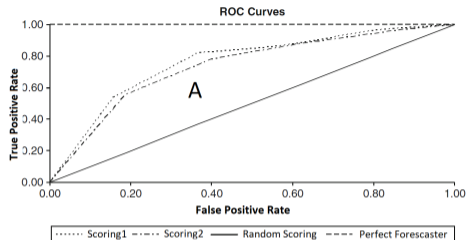
Confusion matrix (in R packages, it is transposed)

		Predicted condition			
		Positive (PP)	Negative (PN)		
Actual condition	Total population $= P + N$			Informedness, bookmaker informedness (BM) $= TPR + TNR - 1$	Prevalence threshold (PT) $= \frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$
	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation	True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, <i>power</i> $= \frac{TP}{P} = 1 - FNR$	False negative rate (FNR), miss rate $= \frac{FN}{P} = 1 - TPR$
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection	False positive rate (FPR), probability of false alarm, fall-out $= \frac{FP}{N} = 1 - TNR$	True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$
	Prevalence $= \frac{P}{P + N}$	Positive predictive value (PPV), <i>precision</i> $= \frac{TP}{PP} = 1 - FDR$	False omission rate (FOR) $= \frac{FN}{PN} = 1 - NPV$	Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$	Negative likelihood ratio (LR-) $= \frac{FNR}{TNR}$
	Accuracy (ACC) $= \frac{TP + TN}{P + N}$	False discovery rate (FDR) $= \frac{FP}{PP} = 1 - PPV$	Negative predictive value (NPV) $= \frac{TN}{PN} = 1 - FOR$	Markedness (MK), deltaP (Δp) $= PPV + NPV - 1$	Diagnostic odds ratio (DOR) = $\frac{LR+}{LR-}$
	Balanced accuracy (BA) = $\frac{TPR + TNR}{2}$	F₁ score $= \frac{2PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$	Fowlkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$	Matthews correlation coefficient (MCC) $= \frac{\sqrt{TPR \times TNR \times PPV \times NPV} - \sqrt{FNR \times FPR \times FOR \times FDR}}$	Threat score (TS), critical success index (CSI), Jaccard index = $\frac{TP}{TP + FN + FP}$

Metrics computed on a test set are intended to estimate some parameter over the general distribution.

- $X = (W, C) \sim F$, i.e., F is the (unknown) multivariate distribution of predictive features and class
- Accuracy ACC of a classifier y_{θ}^+ is a point estimate of $E_F[\mathbb{1}_{y_{\theta}^+(W)=C}] = P_F(y_{\theta}^+(W) = C)$

Probabilistic binary classifier performance metrics



- Binary classifier score $s_\theta(w) \in [0, 1]$ where $s_\theta(w)$ estimates $\eta(w) = P_{\theta_{TRUE}}(C = 1 | W = w)$
- ROC Curve *[Cfr. also Lesson 16]*
 - ▶ $TPR(p) = P(s_\theta(w) \geq p | C = 1)$ and $FPR(p) = P(s_\theta(w) | C = 0)$
 - ▶ ROC Curve is the scatter plot $TPR(p)$ over $FPR(p)$ for p ranging from 1 down to 0
 - ▶ AUC-ROC is the area below the curve What does AUC-ROC estimate?
- Squared error loss or L_2 loss or Brier score: $\frac{1}{n} \sum_i (s_\theta(w_i) - c_i)^2$
- Classifier is calibrated if $P(C = 1 | s_\theta(w) = p) = p$ [classifier-calibration.github.io](https://github.com/classifier-calibration)
 - ▶ Binary Expected Calibration Error (binary-ECE): $\sum_b \frac{|B_b|}{n} |Y_b - S_b|$
 - B_b is the set of i 's in the b^{th} bin, $Y_b = |\{i | i \in B_b, c_i = 1\}| / |B_b|$, $S_b = (\sum_{i \in B_b} s_\theta(w_i)) / |B_b|$