Evaluating Classifier Performance

Tilmann Gneiting

Karlsruhe Institute of Technology (KIT) and Heidelberg Institute for Theoretical Studies (HITS)

> Peter Vogel CSL Behring, Marburg

GridKa School 2019 — The Art of Data 26 August 2019



Heidelberg Institute for Theoretical Studies





(日) (四) (日) (日) (日)

Evaluating classifier performance

- 1. Predicting a binary event
- 2. Receiver operating characteristic (ROC) curves: What are they, and what are they good for?
- 3. Tools for evaluating probabilistic classifiers: A triptych

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Evaluating classifier performance

- 1. Predicting a binary event
- 2. Receiver operating characteristic (ROC) curves: What are they, and what are they good for?
- 3. Tools for evaluating probabilistic classifiers: A triptych

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as

- precipitation tomorrow
- recession
- Champions League final

- medical diagnosis
- recidivism
- 2020 presidential election

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as



- recession
- Champions League final

- medical diagnosis
- recidivism
- 2020 presidential election

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

we let Y = 1 denote a positive and Y = 0 a negative outcome

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as



- recession
- Champions League final

- recidivism
- 2020 presidential election

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

we let Y = 1 denote a positive and Y = 0 a negative outcome

a forecast for Y might take any of the following forms:

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as

- precipitation tomorrow medical diagnosis
- recession
- Champions League final

- recidivism
- 2020 presidential election

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

we let Y = 1 denote a positive and Y = 0 a negative outcome

a forecast for Y might take any of the following forms:

▶ a hard classifier predicts a positive (1) or negative (0) outcome

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as

- precipitation tomorrow
- recession
- Champions League final

- medical diagnosis
- recidivism
- 2020 presidential election

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

we let Y = 1 denote a positive and Y = 0 a negative outcome

- a forecast for Y might take any of the following forms:
 - ▶ a hard classifier predicts a positive (1) or negative (0) outcome
 - ► a probabilistic classifier specifies a probability p ∈ [0, 1] for a positive outcome

consider the (presumably) simplest task in forecasting, namely, the prediction of a binary event, such as

- precipitation tomorrow
- recession
- Champions League final

- medical diagnosis
- recidivism
- 2020 presidential election

we let Y = 1 denote a positive and Y = 0 a negative outcome

- a forecast for Y might take any of the following forms:
 - ▶ a hard classifier predicts a positive (1) or negative (0) outcome
 - ► a probabilistic classifier specifies a probability p ∈ [0, 1] for a positive outcome
 - a feature provides a value x ∈ ℝ, with the understanding that the larger the feature, the more likely a positive outcome

<□ > < □ > < □ > < Ξ > < Ξ > < Ξ > Ξ · のQ@

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

key example in this presentation

key example in this presentation

- 24-hour ahead precipitation forecasts over northern tropical Africa (Vogel 1, 2010)
 - et al. 2018)



▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

key example in this presentation

24-hour ahead precipitation forecasts over northern tropical Africa (Vogel et al. 2018)



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

 based on the numerical weather prediction (NWP) ensemble operated by the European Centre for Medium-Range Weather Forecasts (ECMWF)

key example in this presentation

24-hour ahead precipitation forecasts over northern tropical Africa (Vogel et al. 2018)



based on the numerical weather prediction (NWP) ensemble operated by the European Centre for Medium-Range Weather Forecasts (ECMWF)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

the ensemble has 50 members, each of which provides a hard classifier

key example in this presentation

24-hour ahead precipitation forecasts over northern tropical Africa (Vogel et al. 2018)



based on the numerical weather prediction (NWP) ensemble operated by the European Centre for Medium-Range Weather Forecasts (ECMWF)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

- the ensemble has 50 members, each of which provides a hard classifier
- the resulting probability of precipitation (PoP) forecast equals the fraction of members that predict precipitation

key example in this presentation

24-hour ahead precipitation forecasts over northern tropical Africa (Vogel et al. 2018)



- based on the numerical weather prediction (NWP) ensemble operated by the European Centre for Medium-Range Weather Forecasts (ECMWF)
- ▶ the ensemble has 50 members, each of which provides a hard classifier
- the resulting probability of precipitation (PoP) forecast equals the fraction of members that predict precipitation
- our interest is in comparing the ECMWF PoP forecast to statistically postprocessed and climatological (i.e., purely observation based) forecasts

key example in this presentation

24-hour ahead precipitation forecasts over northern tropical Africa (Vogel et al. 2018)



- based on the numerical weather prediction (NWP) ensemble operated by the European Centre for Medium-Range Weather Forecasts (ECMWF)
- the ensemble has 50 members, each of which provides a hard classifier
- the resulting probability of precipitation (PoP) forecast equals the fraction of members that predict precipitation
- our interest is in comparing the ECMWF PoP forecast to statistically postprocessed and climatological (i.e., purely observation based) forecasts

Vogel, P., Knippertz, P., Fink, A. H., Schlueter, A. and Gneiting, T. (2018). Skill of global raw and postprocessed ensemble predictions of rainfall over northern tropical Africa. *Weather and Forecasting*, 33, 369–388.

Evaluating classifier performance

- 1. Predicting a binary event
- 2. Receiver operating characteristic (ROC) curves: What are they, and what are they good for?
- 3. Tools for evaluating probabilistic classifiers: A triptych

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

▲□▶▲圖▶▲≣▶▲≣▶ ≣ の�?

receiver operating characteristic (ROC) curves are outrageously popular tools for evaluating features (in general) and probabilistic classifiers (in particular)

receiver operating characteristic (ROC) curves are outrageously popular tools for evaluating features (in general) and probabilistic classifiers (in particular)



results from a Web of Science topic search for ''receiver operating characteristic'' or ''ROC''

receiver operating characteristic (ROC) curves are outrageously popular tools for evaluating features (in general) and probabilistic classifiers (in particular)



results from a Web of Science topic search for ''receiver operating characteristic'' or ''ROC''

proposed in signal detection (Egan, Greenberg and Schulman 1961) and cognitive psychology (Swets 1973)

(日) (四) (日) (日) (日)

receiver operating characteristic (ROC) curves are outrageously popular tools for evaluating features (in general) and probabilistic classifiers (in particular)



results from a Web of Science topic search for ''receiver operating characteristic'' or ''ROC''

proposed in signal detection (Egan, Greenberg and Schulman 1961) and cognitive psychology (Swets 1973)

vastly popular ever since Swets (Science, 1988),

receiver operating characteristic (ROC) curves are outrageously popular tools for evaluating features (in general) and probabilistic classifiers (in particular)



results from a Web of Science topic search for ''receiver operating characteristic'' or ''ROC''

proposed in signal detection (Egan, Greenberg and Schulman 1961) and cognitive psychology (Swets 1973)

vastly popular ever since Swets (*Science*, 1988), particularly in the life sciences and machine learning

(日) (四) (日) (日) (日)

▲□▶▲圖▶▲≣▶▲≣▶ ≣ の�?

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

• if X > x we predict a positive outcome (Y = 1)

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- if X > x we predict a positive outcome (Y = 1)
- if $X \le x$ we predict a negative outcome (Y = 0)

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- if X > x we predict a positive outcome (Y = 1)
- if $X \le x$ we predict a negative outcome (Y = 0)

in practice,

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

- if X > x we predict a positive outcome (Y = 1)
- if $X \leq x$ we predict a negative outcome (Y = 0)

in practice,

we are given a dataset

$$\{(x_i, y_i): i=1,\ldots,n\},\$$

where $x_i \in \mathbb{R}$ is the feature and $y_i \in \{0, 1\}$ the associated binary outcome

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

- if X > x we predict a positive outcome (Y = 1)
- if $X \leq x$ we predict a negative outcome (Y = 0)

in practice,

we are given a dataset

$$\{(x_i, y_i): i=1,\ldots,n\},\$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

where $x_i \in \mathbb{R}$ is the feature and $y_i \in \{0, 1\}$ the associated binary outcome

• by thresholding at $x \in \mathbb{R}$, we get a hard classifier

generally, we can use any threshold value x to construct a hard classifier for the binary event Y from the feature X

- if X > x we predict a positive outcome (Y = 1)
- if $X \leq x$ we predict a negative outcome (Y = 0)

in practice,

we are given a dataset

$$\{(x_i, y_i): i = 1, \ldots, n\},\$$

where $x_i \in \mathbb{R}$ is the feature and $y_i \in \{0, 1\}$ the associated binary outcome

- by thresholding at $x \in \mathbb{R}$, we get a hard classifier
- ▶ with true positive (TP : x_i > x, y_i = 1), true negative (TN : x_i < x, y_i = 0), false positive (FP : x_i > x, y_i = 0) and false negative (FN : x_i < x, y_i = 1) instances in the dataset

▲□▶▲圖▶▲≣▶▲≣▶ ≣ の�?

let P and N denote the overall numbers of positive and negative cases in the dataset

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

let P and N denote the overall numbers of positive and negative cases in the dataset

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

given any threshold value x,
let P and N denote the overall numbers of positive and negative cases in the dataset

given any threshold value x,

let TP(x), FN(x), FP(x) and TN(x) denote the respective numbers of true positive, false negative, false positive and true negative cases

let P and N denote the overall numbers of positive and negative cases in the dataset

given any threshold value x,

- let TP(x), FN(x), FP(x) and TN(x) denote the respective numbers of true positive, false negative, false positive and true negative cases
- the hit rate (HR) at x is

$$HR(x) = \frac{TP(x)}{TP(x) + FN(x)} = \frac{TP(x)}{P}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

let P and N denote the overall numbers of positive and negative cases in the dataset

given any threshold value x,

- let TP(x), FN(x), FP(x) and TN(x) denote the respective numbers of true positive, false negative, false positive and true negative cases
- the hit rate (HR) at x is

$$\mathsf{HR}(x) = \frac{\mathsf{TP}(x)}{\mathsf{TP}(x) + \mathsf{FN}(x)} = \frac{\mathsf{TP}(x)}{\mathsf{P}}$$

the false alarm rate (FAR) at x is

$$FAR(x) = \frac{FP(x)}{FP(x) + TN(x)} = \frac{FP(x)}{N}$$

let P and N denote the overall numbers of positive and negative cases in the dataset

given any threshold value x,

- let TP(x), FN(x), FP(x) and TN(x) denote the respective numbers of true positive, false negative, false positive and true negative cases
- the hit rate (HR) at x is

$$\mathsf{HR}(x) = \frac{\mathsf{TP}(x)}{\mathsf{TP}(x) + \mathsf{FN}(x)} = \frac{\mathsf{TP}(x)}{\mathsf{P}}$$

the false alarm rate (FAR) at x is

$$FAR(x) = \frac{FP(x)}{FP(x) + TN(x)} = \frac{FP(x)}{N}$$

the ROC curve plots the hit rate HR(x) vs. the false alarm rate FAR(x) as the decision threshold x varies

・ロト ・西ト ・ヨト ・ヨー うへぐ

let P and N denote the overall numbers of positive and negative cases in the dataset

given any threshold value x,

- let TP(x), FN(x), FP(x) and TN(x) denote the respective numbers of true positive, false negative, false positive and true negative cases
- the hit rate (HR) at x is

$$\mathsf{HR}(x) = \frac{\mathsf{TP}(x)}{\mathsf{TP}(x) + \mathsf{FN}(x)} = \frac{\mathsf{TP}(x)}{\mathsf{P}}$$

the false alarm rate (FAR) at x is

$$FAR(x) = \frac{FP(x)}{FP(x) + TN(x)} = \frac{FP(x)}{N}$$

the ROC curve plots the hit rate HR(x) vs. the false alarm rate FAR(x) as the decision threshold x varies

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

appealing interpretation as the probability that a (randomly chosen) feature value under a positive outcome is larger than a (randomly chosen) value under a negative outcome:

$$AUC = \mathbb{P}(X' > X \mid Y' = 1, Y = 0)$$

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

appealing interpretation as the probability that a (randomly chosen) feature value under a positive outcome is larger than a (randomly chosen) value under a negative outcome:

$$\mathsf{AUC} = \mathbb{P}(X' > X \mid Y' = 1, Y = 0)$$

• AUC = (D + 1)/2 in terms of Somers' D (Somers 1962)

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

appealing interpretation as the probability that a (randomly chosen) feature value under a positive outcome is larger than a (randomly chosen) value under a negative outcome:

$$\mathsf{AUC} = \mathbb{P}(X' > X \mid Y' = 1, Y = 0)$$

- AUC = (D + 1)/2 in terms of Somers' D (Somers 1962)
- AUC = 1/2 for a useless feature that is independent of the binary outcome

the area under the ROC curve (AUC) is a positively oriented measure of predictive ability

appealing interpretation as the probability that a (randomly chosen) feature value under a positive outcome is larger than a (randomly chosen) value under a negative outcome:

$$\mathsf{AUC} = \mathbb{P}(X' > X \mid Y' = 1, Y = 0)$$

- AUC = (D + 1)/2 in terms of Somers' D (Somers 1962)
- AUC = 1/2 for a useless feature that is independent of the binary outcome

<ロト < 個 ト < 臣 ト < 臣 ト 三 の < @</p>

ROC curve and AUC for 24-hour probability of precipitation (PoP) forecasts from the ECMWF ensemble over West Sahel

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

ROC curve and AUC for 24-hour probability of precipitation (PoP) forecasts from the ECMWF ensemble over West Sahel

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Vogel et al. (2018)



ROC curve and AUC for 24-hour probability of precipitation (PoP) forecasts from the ECMWF ensemble over West Sahel

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @

Vogel et al. (2018)



ROC curve and AUC for 24-hour probability of precipitation (PoP) forecasts from the ECMWF ensemble over West Sahel

Vogel et al. (2018)



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○三 の々で

ROC curve and AUC for 24-hour probability of precipitation (PoP) forecasts from the ECMWF ensemble over West Sahel



◆□ > ◆□ > ◆豆 > ◆豆 > ̄豆 = のへで

▲□▶▲圖▶▲≣▶▲≣▶ ≣ の�?

formal setting

- X real-valued feature
- Y binary outcome
- \mathbb{P} joint distribution of (X, Y)

 $F_1(x) = \mathbb{P}(X \le x \mid Y = 1)$ $F_0(x) = \mathbb{P}(X \le x \mid Y = 0)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

formal setting

 $\begin{array}{ll} X \ \text{real-valued feature} \\ Y \ \text{binary outcome} \\ \mathbb{P} \ \text{joint distribution of } (X, Y) \end{array} \end{array} \begin{array}{ll} F_1(x) = \mathbb{P}(X \leq x \mid Y = 1) \\ F_0(x) = \mathbb{P}(X \leq x \mid Y = 0) \end{array}$

the raw ROC diagnostic is the set of all points of the form $(FAR(x), HR(x)) \in [0, 1] \times [0, 1]$ threshold $x \in \mathbb{R}$, $FAR(x) = 1 - F_0(x)$, $HR(x) = 1 - F_1(x)$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

formal setting

 $\begin{array}{ll} X \ \text{real-valued feature} \\ Y \ \text{binary outcome} \\ \mathbb{P} \ \text{joint distribution of } (X, Y) \end{array} \end{array} \begin{array}{ll} F_1(x) = \mathbb{P}(X \leq x \mid Y = 1) \\ F_0(x) = \mathbb{P}(X \leq x \mid Y = 0) \end{array}$

the raw ROC diagnostic is the set of all points of the form $(FAR(x), HR(x)) \in [0, 1] \times [0, 1]$ threshold $x \in \mathbb{R}$, $FAR(x) = 1 - F_0(x)$, $HR(x) = 1 - F_1(x)$

(1100) (110) (110) (110) (110) (110) (110)

the ROC curve is the linearly interpolated raw ROC diagnostic



interpretation as function: for continuous, strictly increasing F_0 and F_1 ,

 $R(\alpha) = 1 - F_1(F_0^{-1}(1 - \alpha)), \quad \alpha = FAR(x) \in (0, 1)$



ROC curve



interpretation as function: for continuous, strictly increasing F_0 and F_1 ,

$$R(\alpha) = 1 - F_1(F_0^{-1}(1 - \alpha)), \quad \alpha = FAR(x) \in (0, 1)$$

ensuing math fact: characterization of ROC curves



ROC curve



interpretation as function: for continuous, strictly increasing F_0 and F_1 ,

$$R(\alpha) = 1 - F_1(F_0^{-1}(1 - \alpha)), \quad \alpha = FAR(x) \in (0, 1)$$

ensuing math fact: characterization of ROC curves

invariance of ROC curves and AUC under

- changes in class proportions
- strictly increasing transformations of the feature X



ROC curve



interpretation as function: for continuous, strictly increasing F_0 and F_1 ,

$$R(\alpha) = 1 - F_1(F_0^{-1}(1 - \alpha)), \quad \alpha = FAR(x) \in (0, 1)$$

ensuing math fact: characterization of ROC curves

invariance of ROC curves and AUC under

- changes in class proportions
- strictly increasing transformations of the feature X

consequence: ROC curves and AUC

- apply to real-valued or ordinal features X on arbitrary scales, but
- do not consider calibration nor economic value for probabilistic classifiers







<ロ> < 団> < 団> < 豆> < 豆> < 豆> < 豆> < </p>

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ



24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●



24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ ▲ 三 ● ● ●



competing probability of precipitation (PoP) forecasts

 ENS ECMWF NWP ensemble

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @



- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics
- BMA postprocessed by Bayesian model averaging

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで



- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics
- BMA postprocessed by Bayesian model averaging
- EPC extended probabilistic climatology

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)



competing probability of precipitation (PoP) forecasts

- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics
- BMA postprocessed by Bayesian model averaging
- EPC extended probabilistic climatology



・ロト ・ 同ト ・ ヨト ・ ヨト

э

Evaluating classifier performance

- 1. Predicting a binary event
- 2. Receiver operating characteristic (ROC) curves: What are they, and what are they good for?
- 3. Tools for evaluating probabilistic classifiers: A triptych

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

ROC curves address potential predictive ability (only)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

invariance under strictly monotone transformations has stark implications:
invariance under strictly monotone transformations has stark implications:

 ROC curves and AUC can be used to assess the potential predictive ability of any real-valued feature

invariance under strictly monotone transformations has stark implications:

- ROC curves and AUC can be used to assess the potential predictive ability of any real-valued feature
- however, in the case of probabilistic classifiers, both calibration and actual value are ignored

invariance under strictly monotone transformations has stark implications:

- ROC curves and AUC can be used to assess the potential predictive ability of any real-valued feature
- however, in the case of probabilistic classifiers, both calibration and actual value are ignored

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

hence, ROC curves should be used in concert with reliability diagrams and Murphy diagrams

invariance under strictly monotone transformations has stark implications:

- ROC curves and AUC can be used to assess the potential predictive ability of any real-valued feature
- however, in the case of probabilistic classifiers, both calibration and actual value are ignored
- hence, ROC curves should be used in concert with reliability diagrams and Murphy diagrams

reliability diagrams

assess calibration, by plotting the empirical conditional event frequency vs. the predicted probability

invariance under strictly monotone transformations has stark implications:

- ROC curves and AUC can be used to assess the potential predictive ability of any real-valued feature
- however, in the case of probabilistic classifiers, both calibration and actual value are ignored
- hence, ROC curves should be used in concert with reliability diagrams and Murphy diagrams

reliability diagrams

assess calibration, by plotting the empirical conditional event frequency vs. the predicted probability

Murphy diagrams

assess actual economic value, by considering all (!) proper scoring rules simultaneously

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

we take scores to be negatively oriented: the smaller, the better

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

we take scores to be negatively oriented: the smaller, the better

a proper scoring rule rewards honest and careful forecasting: truth telling is the best strategy in expectation

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

we take scores to be negatively oriented: the smaller, the better

a proper scoring rule rewards honest and careful forecasting: truth telling is the best strategy in expectation

every proper scoring rule can be represented as a weighted mixture over elementary scores

$$\mathsf{S}_{\theta}(\boldsymbol{p},\boldsymbol{y}) = \begin{cases} \theta, & \boldsymbol{y} = \boldsymbol{0}, \ \boldsymbol{p} > \theta, \\ 1 - \theta, & \boldsymbol{y} = \boldsymbol{1}, \ \boldsymbol{p} \leq \theta, \\ 0, & \text{otherwise}, \end{cases}$$

where the index $\theta \in (0, 1)$ reflects a decision maker's cost ratio; e.g.

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

we take scores to be negatively oriented: the smaller, the better

a proper scoring rule rewards honest and careful forecasting: truth telling is the best strategy in expectation

every proper scoring rule can be represented as a weighted mixture over elementary scores

$$\mathsf{S}_{\theta}(\boldsymbol{p},\boldsymbol{y}) = \begin{cases} \theta, & \boldsymbol{y} = \boldsymbol{0}, \ \boldsymbol{p} > \theta, \\ 1 - \theta, & \boldsymbol{y} = \boldsymbol{1}, \ \boldsymbol{p} \leq \theta, \\ 0, & \text{otherwise}, \end{cases}$$

where the index $\theta \in (0, 1)$ reflects a decision maker's cost ratio; e.g.

For the quadratic or Brier score, S(p, y) = (p − y)², the mixture is (twice) uniform over θ ∈ (0, 1)

a scoring rule is a function S(p, y) that assigns a numerical score to a probability forecast p with binary outcome y

we take scores to be negatively oriented: the smaller, the better

a proper scoring rule rewards honest and careful forecasting: truth telling is the best strategy in expectation

every proper scoring rule can be represented as a weighted mixture over elementary scores

$$\mathsf{S}_{ heta}(p,y) = egin{cases} heta, & y = 0, \ p > heta, \ 1 - heta, & y = 1, \ p \leq heta, \ 0, & ext{otherwise}, \end{cases}$$

where the index $\theta \in (0, 1)$ reflects a decision maker's cost ratio; e.g.

For the quadratic or Brier score, S(p, y) = (p − y)², the mixture is (twice) uniform over θ ∈ (0, 1)

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

For the logarithmic score, S(p, y) = − log p, the mixture measure has density ∝ (θ(1 − θ))⁻¹

a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

the Brier score equals (twice) the area under the Murphy curve

a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

- the Brier score equals (twice) the area under the Murphy curve
- ▶ the widely reported misclassification rate equals (twice) the height of the Murphy curve at $\theta = 1/2$

a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

- the Brier score equals (twice) the area under the Murphy curve
- ▶ the widely reported misclassification rate equals (twice) the height of the Murphy curve at $\theta = 1/2$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ● ●

a Murphy diagram plots the curves of competing forecasters in a single graph, as implemented in the R package murphydiagram

a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

- the Brier score equals (twice) the area under the Murphy curve
- ► the widely reported misclassification rate equals (twice) the height of the Murphy curve at $\theta = 1/2$

a Murphy diagram plots the curves of competing forecasters in a single graph, as implemented in the R package murphydiagram



a Murphy curve plots the mean elementary score S_{θ} of a probabilistic classifier as a function of $\theta \in (0, 1)$

- the Brier score equals (twice) the area under the Murphy curve
- the widely reported misclassification rate equals (twice) the height of the Murphy curve at θ = 1/2

a Murphy diagram plots the curves of competing forecasters in a single graph, as implemented in the R package murphydiagram



covers all economic scenarios simultaneously and eliminates the need to choose a proper scoring rule (Murphy 1977; Ehm et al. 2016)

ROC curve, reliability diagram, and Murphy diagram ... a triptych

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへぐ

ROC curve, reliability diagram, and Murphy diagram ...a triptych according to Wikipedia, a triptych is "a piece of art [...] that is divided into three sections"

ROC curve, reliability diagram, and Murphy diagram

according to Wikipedia, a triptych is "a piece of art $[\dots]$ that is divided into three sections"



Merode Altarpiece by Robert Campin, ca. 1427–32 Metropolitan Museum of Art

Source: Wikimedia https://commons.wikimedia.org/wiki/File:Robert_Campin_-_L%27_Annonciation_-_1425.jpg

Example

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)

▲ロ ▶ ▲周 ▶ ▲ 国 ▶ ▲ 国 ▶ ● の Q @



competing probability of precipitation (PoP) forecasts

- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics
- BMA postprocessed by Bayesian model averaging
- EPC extended probabilistic climatology

Example

24-hour precipitation forecasts over the West Sahel region in northern tropical Africa in monsoon season 2014 (Vogel et al. 2018)



competing probability of precipitation (PoP) forecasts

- ENS ECMWF NWP ensemble
- EMOS postprocessed by ensemble model output statistics
- BMA postprocessed by Bayesian model averaging
- EPC extended probabilistic climatology



・ロト ・ 国 ト ・ ヨ ト ・ ヨ ト

э

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

A D > A P > A B > A B >

э

ROC curves assess potential predictive ability



ROC curves assess potential predictive ability

Murphy diagrams visualize actual (normalized) costs for a decision maker with expense ratio $\theta/(1-\theta)$





イロト 不得 トイヨト イヨト

э

ROC curves assess potential predictive ability

reliability diagrams demonstrate calibration

Murphy diagrams visualize actual (normalized) costs for a decision maker with expense ratio $\theta/(1-\theta)$



▲□▶ ▲□▶ ▲臣▶ ★臣▶ = 臣 = のへで

Selected references and code

Gneiting, T. and Vogel, P. (2018). Receiver operating characteristic (ROC) curves. Preprint, arXiv:1809.04808.

Ehm, W., Gneiting, T., Jordan, A. and Krüger, F. (2016), Of quantiles and expectiles: Consistent scoring functions, Choquet representations and forecast rankings (with discussion and rejoinder), *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78, 505–562.

Vogel, P., Knippertz, P., Fink, A. H., Schlueter, A. and Gneiting, T. (2018). Skill of global raw and postprocessed ensemble predictions of rainfall over northern tropical Africa. *Weather and Forecasting*, 33, 369–388.

some useful R packages at CRAN: pROC, ROCR, verification, scoringRules, murphydiagram