

# From Uncertainty to Precision: Enhancing Binary Classifier Performance through Calibration.

A. Fernandes Machado<sup>1</sup>   A. Charpentier<sup>1</sup>   E. Flachaire<sup>2</sup>   E. Gallic<sup>2</sup>   F. Hu<sup>3</sup>

<sup>1</sup>Université du Québec à Montréal

<sup>2</sup>Aix-Marseille School of Economics, Aix-Marseille Univ.

<sup>3</sup>Milliman France



58th Annual Canadian Economics Association Meetings, May 31st, 2024

# Calibration: intuition (1/2)

*"yPCqCS - {CE <P- ^<C  
bHq S^ zb\ bqp.i"*  
(Dawid, 1982)

# Calibration: intuition (1/2)

“ $y_{PGCS} - \{CE \langle P \wedge C$   
 $bHq S^z b \backslash b q p . i$ ”  
 (Dawid, 1982)











Thu 30	66°/51°	 Sunny	2%	 N 11 mph	▼
Fri 31	71°/54°	 Sunny	1%	 SW 12 mph	▼
Sat 01	72°/61°	 Mostly Cloudy	3%	 S 10 mph	▼
Sun 02	69°/61°	 AM Showers	49%	 E 10 mph	▼
Mon 03	70°/64°	 Partly Cloudy	16%	 E 9 mph	▼

Figure 1: Weather Forecasts on Tuesday, March 2024. Source: The Weather Channel.

# Calibration: intuition (1/2)

“ $y_{PGCS} - \{ \mathcal{C} \} \langle P \rangle \wedge \mathcal{C}$   
 $bHq S^z b \backslash b q p . i$ ”  
 (Dawid, 1982)











Thu 30	66°/51°	 Sunny	2%	 N 11 mph	▼
Fri 31	71°/54°	 Sunny	1%	 SW 12 mph	▼
Sat 01	72°/61°	 Mostly Cloudy	3%	 S 10 mph	▼
Sun 02	69°/61°	 AM Showers	49%	 E 10 mph	▼
Mon 03	70°/64°	 Partly Cloudy	16%	 E 9 mph	▼

Figure 1: Weather Forecasts on Tuesday, March 2024. Source: The Weather Channel.

Consider a sequence of weather forecasts  $\hat{s}(\ddagger_t)$ , where  $t = 1, \dots, T$  denotes the days of forecast and  $\ddagger$  represents characteristics used in forecasting.

## Calibration: intuition (2/2)

Within this sequence, we focus on days where  $\hat{s}(\dagger_i)$  closely approximates 30%.

By assuming an infinite sequence, we can determine the long-term proportion  $p$  of days where the forecasted event actually occurred.



## Calibration: intuition (2/2)

Within this sequence, we focus on days where  $\hat{s}(\dagger_i)$  closely approximates 30%.

By assuming an infinite sequence, we can determine the long-term proportion  $p$  of days where the forecasted event actually occurred.



## Calibration: intuition (2/2)

Within this sequence, we focus on days where  $\hat{s}(\dagger_i)$  closely approximates 30%.

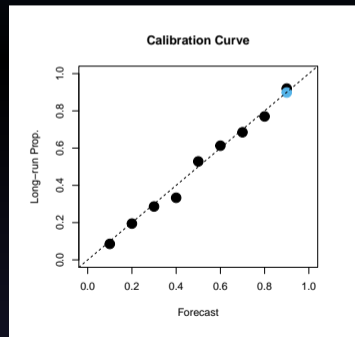
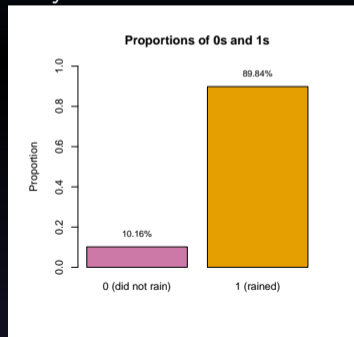
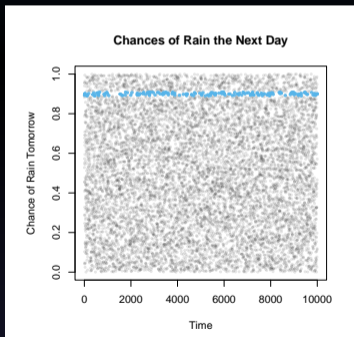
By assuming an infinite sequence, we can determine the long-term proportion  $p$  of days where the forecasted event actually occurred.



# Calibration: intuition (2/2)

Within this sequence, we focus on days where  $\hat{s}(\dagger_i)$  closely approximates 30%.

By assuming an infinite sequence, we can determine the long-term proportion  $p$  of days where the forecasted event actually occurred.





# Motivations

- We are interested in being able to  $\mathbb{P}(S < z | C)$  between rainy/not rainy days.

# Motivations

- We are interested in being able to  $\mathbb{P}(S < z) \approx \mathbb{P}(S - z < 0)$  between rainy/not rainy days.
- We are also interested in the  $\sim^{\wedge} \mathbb{P}(C \leq L) \approx \mathbb{P}(W)$  Other examples include:

# Motivations

- We are interested in being able to  $\mathbb{P}(S=1|C)$  between rainy/not rainy days.
- We are also interested in the  $\mathbb{P}(S=1|C)$  Other examples include:
  - does this patient have a disease or not (Van Calster et al., 2019)?

# Motivations

- We are interested in being able to  $\mathbb{P}(S=1|Z)$  between rainy/not rainy days.
- We are also interested in the  $\mathbb{P}(S=1|Z)$  Other examples include:
  - does this patient have a disease or not (Van Calster et al., 2019)?
  - will this insured have an accident within the next year?

# Motivations

- We are interested in being able to  $\mathbb{P}(S=1|C)$  between rainy/not rainy days.
- We are also interested in the  $\mathbb{P}(S=1|C)$  Other examples include:
  - does this patient have a disease or not (Van Calster et al., 2019)?
  - will this insured have an accident within the next year?
  - what is the probability for this individual to receive the treatment/control?

# Motivations

- We are interested in being able to  $\mathbb{P}(S=1|C)$  between rainy/not rainy days.
- We are also interested in the  $\mathbb{P}(S=1|C)$  Other examples include:
  - does this patient have a disease or not (Van Calster et al., 2019)?
  - will this insured have an accident within the next year?
  - what is the probability for this individual to receive the treatment/control?
- In such cases, it is important that the  $\mathbb{P}(S=1|C)$  can be interpreted as  $\mathbb{P}(S=1|C)$ .

# Motivations

- We are interested in being able to  $\mathbb{P}(S=1|C)$  between rainy/not rainy days.
- We are also interested in the  $\mathbb{P}(S=1|C)$  Other examples include:
  - does this patient have a disease or not (Van Calster et al., 2019)?
  - will this insured have an accident within the next year?
  - what is the probability for this individual to receive the treatment/control?
- In such cases, it is important that the  $\mathbb{P}(S=1|C)$  can be interpreted as  $\mathbb{P}(S=1|C)$ .
- This might become a problem when using  $\mathbb{P}(S=1|C)$  based on  $\mathbb{P}(S=1|C)$ .

# This Talk

What the remainder of the talk is about:

- Reviewing of ways to  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$  for a  $\mathcal{L}(y, \hat{y}) = \mathbb{1}_{y \neq \hat{y}}$



# This Talk

What the remainder of the talk is about:

- Reviewing of ways to  $\int_{\mathcal{C}} \mathbb{1}_{\{y \neq \hat{y}\}} p(y) dy$  for a  $\mathcal{C}$ -valued  $p$
- Proposing a new metric based on  $\int_{\mathcal{C}} \mathbb{1}_{\{y \neq \hat{y}\}} p(y) dy$ : the Local Calibration Score.

# This Talk

What the remainder of the talk is about:

- Reviewing of ways to  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$  for a  $\mathcal{L}(y, \hat{y}) = \mathbb{1}_{y \neq \hat{y}}$
- Proposing a new metric based on  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$ : the Local Calibration Score.
- Observing the  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$  on standard performance metrics.

# This Talk

What the remainder of the talk is about:

- Reviewing of ways to  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$  for a  $\mathcal{L}(y, \hat{y}) = \mathbb{1}_{y \neq \hat{y}}$
- Proposing a new metric based on  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$ : the Local Calibration Score.
- Observing the  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$  on standard performance metrics.
- Examining calibration for  $\int_{\mathcal{C}} \mathcal{L}(y, \hat{y}) p(y) dy$ .

# Take away results

- Our new metric, the  $\mathbb{E}[\max_{c \in \mathcal{C}} |p_c - Y_c|]$  offers a more flexible way to visualise and measure calibration than methods based on empirical quantiles.
- $\mathbb{E}[\max_{c \in \mathcal{C}} |p_c - Y_c|]$ : when training classifiers, looking at calibration of models should not be disregarded.

# Roadmap

- 1 Introduction
- 2 Calibration
  - Definition
  - Measuring Calibration
- 3 Impact of Poor Calibration
- 4 Calibration and Tree-Based Methods

# Roadmap

## Calibration

# Setup

- Let us consider a  $\mathcal{D}$  whose observations are denoted  $d_i = 1$  if the event occurs, and  $d_i = 0$  otherwise, where  $i$  denotes the  $i$ th observations.

# Setup

- Let us consider a  $D$  whose observations are denoted  $d_i = 1$  if the event occurs, and  $d_i = 0$  otherwise, where  $i$  denotes the  $i$ th observations.
- Let us further assume that the probability of the event  $d_i = 1$  depends on

$$p_i = s(\mathbf{x}_i)$$

where, with sample size  $n > 0$ ,  $i = 1, \dots, n$  represents individuals, and  $\mathbf{x}_i$  the characteristics.



## Predicting risks

- To ~~CS~~ - ~~CS~~ ~~CS~~ ~~CS~~ ~~CS~~ ~~CS~~ we can use a statistical model (*GLi*, a GLM) or a machine learning model (*GLi*, a random forest).

## Predicting risks

- To  $\mathbb{P}(Y=1|X=x)$  we can use a statistical model ( $GLM$ , a GLM) or a machine learning model ( $GLM$ , a random forest).
- These models  $\hat{p}_i \in [0, 1]$ , allowing the  $\hat{p}_i$  of observations based on the estimated probability of the event.

## Predicting risks

- To **assess the risk** we can use a statistical model ( $GLM$ , a GLM) or a machine learning model ( $GLM$ , a random forest).
- These models **output**  $\hat{p}_i \in [0, 1]$ , allowing the **classification** of observations based on the estimated probability of the event.
- By setting a probability threshold  $\tau$  in  $[0, 1]$ , one can predict the **class** of each observation: 1 if the event occurs, and 0 otherwise:

$$\hat{d}_i = \begin{cases} 1, & \text{if } \hat{p}_i \geq \tau \\ 0, & \text{if } \hat{p}_i < \tau \end{cases} .$$

## Predicting risks

- To **assess the risk** of an event, we can use a statistical model ( $GLM$ , a GLM) or a machine learning model ( $GLM$ , a random forest).
- These models **output**  $\hat{s}(x_i) \in [0, 1]$ , allowing the **prediction** of observations based on the estimated probability of the event.
- By setting a probability threshold  $\tau$  in  $[0, 1]$ , one can predict the **class** of each observation: 1 if the event occurs, and 0 otherwise:

$$\hat{d}_i = \begin{cases} 1, & \text{if } \hat{s}(x_i) \geq \tau \\ 0, & \text{if } \hat{s}(x_i) < \tau \end{cases} .$$

- However, if the model is not **calibrated**, the **scores** cannot be interpreted as probabilities.

# Definition

## Calibration of a Binary Classifier (Schervish, 1989)

For a binary variable  $D$ , a model is well-calibrated when

$$\mathbb{E}[D \mid \hat{s}(\dagger) = p] = p, \quad \forall p \in [0, 1] . \quad (1)$$

# Definition

## Calibration of a Binary Classifier (Schervish, 1989)

For a binary variable  $D$ , a model is well-calibrated when

$$\mathbb{E}[D \mid \hat{s}(\dagger) = p] = p, \quad \forall p \in [0, 1] . \quad (1)$$

Note: conditioning by  $\{\hat{s}(\dagger) = p\}$  leads to the concept of (local) calibration; however, as discussed by Bai et al., 2021,  $\{\hat{s}(\dagger) = p\}$  is *isi* a null mass event. Thus, calibration should be understood in the sense that

$$\mathbb{E}[D \mid \hat{s}(\dagger) = p] \xrightarrow{a.s.} p \text{ when } n \rightarrow \infty ,$$

meaning that, asymptotically, the model is well-calibrated, or locally well-calibrated in  $p$ , for any  $p$ .

## Visual approach: calibration curve

- Estimation of  $\mathcal{C}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{S}(\cdot) = p]$  on predicted scores  $\hat{S}(\cdot)$ ):

$$\mathcal{C} : \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \mathcal{C}(p) := \mathbb{E}[D | \hat{S}(\cdot) = p] \end{cases} \quad (2)$$

## Visual approach: calibration curve

- Estimation of  $\mathcal{C}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{S}(\mathbf{x}) = p]$  on predicted scores  $\hat{S}(\mathbf{x})$ ):

$$\mathcal{C} : \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \mathcal{C}(p) := \mathbb{E}[D | \hat{S}(\mathbf{x}) = p] \end{cases} \quad (2)$$

- $\mathcal{C}$ ;  $\mathcal{P}$ - $\mathcal{C}$  having enough observations with identical scores is difficult.



## Visual approach: calibration curve

- Estimation of  $\hat{c}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{c}(\cdot) = p]$  on predicted scores  $\hat{c}(\cdot)$ ):

$$\hat{c}: \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \hat{c}(p) := \mathbb{E}[D | \hat{c}(\cdot) = p] \end{cases} \quad (2)$$

- **Problem**: having enough observations with identical scores is difficult.
- **Solution**: grouping obs. into  $B$  bins, defined by the  $B$  quantiles of predicted scores:

## Visual approach: calibration curve

- Estimation of  $\hat{c}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{s}(\cdot) = p]$  on predicted scores  $\hat{s}(\cdot)$ ):

$$\hat{c} : \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \hat{c}(p) := \mathbb{E}[D | \hat{s}(\cdot) = p] \end{cases} \quad (2)$$

- $\hat{c}(\cdot)$ : having enough observations with identical scores is difficult.
- $\hat{c}(\cdot)$ : grouping obs. into  $B$  bins, defined by the  $\hat{s}(\cdot)$  of predicted scores:
  - The average of observed values ( $\bar{d}_b$  with  $b \in \{1, \dots, B\}$ ), in each bin  $b$  can then be compared with the central value of the bin.

## Visual approach: calibration curve

- Estimation of  $\hat{c}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{y} = p]$  on predicted scores  $\hat{y}$ ):

$$\hat{c}: \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \hat{c}(p) := \mathbb{E}[D | \hat{y} = p] \end{cases} \quad (2)$$

- **Problem**: having enough observations with identical scores is difficult.
- **Solution**: grouping obs. into  $B$  bins, defined by the **edges** of predicted scores:
  - The average of observed values ( $\bar{d}_b$  with  $b \in \{1, \dots, B\}$ ), in each bin  $b$  can then be compared with the central value of the bin.
  - **Reliability diagram** (Wilks, 1990): middle of each bin on the x-axis, averages of corresponding observations on the y-axis.

## Visual approach: calibration curve

- Estimation of  $\hat{c}(\cdot)$  (which measures  $\mathbb{E}[D | \hat{y} = p]$  on predicted scores  $\hat{y}$ ):

$$\hat{c}: \begin{cases} [0, 1] \rightarrow [0, 1] \\ p \mapsto \hat{c}(p) := \mathbb{E}[D | \hat{y} = p] \end{cases} \quad (2)$$

- $\hat{c}$ : having enough observations with identical scores is difficult.
- $\hat{c}$ : grouping obs. into  $B$  bins, defined by the  $\hat{y}$  of predicted scores:
  - The average of observed values ( $\bar{d}_b$  with  $b \in \{1, \dots, B\}$ ), in each bin  $b$  can then be compared with the central value of the bin.
  - $\hat{c}$ : reliability diagram (Wilks, 1990): middle of each bin on the x-axis, averages of corresponding observations on the y-axis.
  - When the model is  $\hat{y} = p$ , all  $B$  points lie on the  $\hat{c}$

# Metrics (1/2)

1 t T 2 + i 2 / \* H B # ` i B Q M S ` F Q ` K Q M 1 L \* 1 2 B M B 2 V H X - k y R 8

$$1 * 1 = \sum_{b=1}^B \frac{n_b}{n} \quad | \quad +(\mathbf{b}) - + Q(\mathbf{M}) \uparrow$$

r ? 2 ` n 2 B b i ? 2 b K T H y 2 B b k 2 2 M m K # 2 ` Q 7 Q # b 2 \in \{1, B, Q, M\} b X B M # B M

## Metrics (1/2)

1 t T 2 + i 2 / \* H B # ` i B Q M S ` F Q ` K Q M 1 L \* 1 2 B M B 2 V H X - k y R 8

$$1 * 1 = \sum_{b=1}^B \frac{n_b}{n} | \quad +(\mathbf{b}) - +Q(M) \uparrow$$

r ? 2 ` r 2 B b i ? 2 b K T H y 2 B b B k 2 2 M m K # 2 ` Q 7 Q # b 2 ` e p { i i , B , Q , M } b X B M # B M

, <<~q <%o +(\mathbf{b}), h ? 2 p 2 ` ; 2 Q 7 2 K T B ` B + H

T ` Q # # B H B i B 2 b Q ` 7 ` + T B 2 Q M b i Q / 7 + Q ` ` 2 + i H v

+ H b b 2 b X

$$+(\mathbf{b}) = \frac{1}{n_b} \sum_{i \in \mathcal{I}_b} \mathbb{1}_{\hat{d}_i = d_i} \quad U j v$$

h ? 2 T ` 2 / B + i 2 \hat{d}\_i \uparrow 10 ` b Q # b 2 i p B b B Q M

/ 2 i 2 ` K B M 2 / # b 2 / Q M + H b b B } + i B Q M i ? ` 2 b ? Q H /

$\tau \in [0, 1]$  r ? 2 ` \hat{d}\_i = 1 B \hat{s}(\mathbf{t}\_i) \geq \tau M \emptyset

Q i ? 2 ` r B b 2 X

# Metrics (1/2)

1 t T 2 + i 2 / \* H B # ` i B Q M S ` F Q ` K Q M 1 L \* 1 2 B M B 2 V H X - k y R 8

$$1 * 1 = \sum_{b=1}^B \frac{n_b}{n} | \quad +(\hat{b}) - + Q(M) \uparrow$$

r ? 2 ` r 2 B b i ? 2 b K T H y B b B k 2 2 M m K # 2 ` Q 7 Q # b 2 ` e p { i i , B , Q , M } , B M # B M

, <<~q <%o +(\hat{b}), h ? 2 p 2 ` ; 2 Q 7 2 K T B ` B + H

T ` Q # # B H B i B 2 b Q ` 7 ` + T B Q M b i Q / 7 + Q ` b 2 a i H y + Q ( M ) 7 A M / B + i 2 b i ? 2 K Q / 2 H + H b b 2 b X

$$+(\hat{b}) = \frac{1}{n_b} \sum_{i \in \mathcal{J}_b} \mathbb{1}_{\hat{d}_i = d_i}$$

p 2 ` ; 2 + Q M } / 2 M + 2 b r B v ? B Q ` # ; B M M ; U j V T ` 2 / B + i 2 / b + Q ` 2 b X

h ? 2 T ` 2 / B + i 2 \hat{d}\_i + 7 Q ` b Q # b 2 ` i p B b B Q M / 2 i 2 ` K B M 2 / # b 2 / Q M + H b b B } + i B Q M i ? ` 2 b ? Q H /

$$+ Q ( M ) \uparrow = \frac{1}{n_b} \sum_{i \in \mathcal{J}_b} \hat{s}(t_i)$$

$\tau \in [0, 1]$  r ? 2 ` \hat{d}\_i = 1 B \hat{s}(t\_i) \geq \tau M \emptyset

Q i ? 2 ` r B b 2 X

## Metrics (2/2)

### Brier Score (Brier, 1950)

The **Brier Score** does not depend on bins and is defined as:

$$BS = \frac{1}{n} \sum_{i=1}^n (d_i - \hat{s}(x_i))^2 \quad (4)$$

where  $d_i$  is the observed event and  $\hat{s}(x_i)$  the estimated score.



## Metrics (2/2)

### Brier Score (Brier, 1950)

The **Brier Score** does not depend on bins and is defined as:

$$BS = \frac{1}{n} \sum_{i=1}^n (d_i - \hat{s}(x_i))^2 \quad (4)$$

where  $d_i$  is the observed event and  $\hat{s}(x_i)$  the estimated score.

### Mean Squared Error (MSE)

By substituting the observed event  $d_i$  by the true probability  $p_i$  (which can only be observed in an experimental setup), the metric becomes the MSE:

$$\text{True MSE} = \frac{1}{n} \sum_{i=1}^n (p_i - \hat{s}(x_i))^2 \quad (5)$$

## Our Approach: Smoother Visualization Technique

We propose an alternative approach to visualize model calibration, aiming for a smoother visualization (Loader, 1999).

- Measuring calibration consists in estimating a local regression seems appropriate.

## Our Approach: Smoother Visualization Technique

We propose an alternative approach to visualize model calibration, aiming for a smoother visualization (Loader, 1999).

- Measuring calibration consists in estimating a local regression seems appropriate.
- Local Regression have been disregarded in high dimensions due to poor properties, but it is appropriate, as in this case with only one predictive feature,  $\hat{s}(x) \in [0, 1]$ .

## Our Approach: Smoother Visualization Technique

We propose an alternative approach to visualize model calibration, aiming for a  $\hat{s}(x)$  (Loader, 1999).

- Measuring calibration consists in estimating a  $\hat{s}(x)$ : a local regression seems appropriate.
- Local Regression have been disregarded in high dimensions due to poor properties, but it is  $\hat{s}(x) \in [0, 1]$ , as in this case with only one predictive feature,  $\hat{s}(x) \in [0, 1]$ .
- Given the number of data points, the precision of quantile binning can be suboptimal when determining the appropriate bin count.

## Our Approach: Smoother Visualization Technique

We propose an alternative approach to visualize model calibration, aiming for a smoother visualization (Loader, 1999).

- Measuring calibration consists in estimating a local regression seems appropriate.
- Local Regression have been disregarded in high dimensions due to poor properties, but it is appropriate as in this case with only one predictive feature,  $\hat{s}(x) \in [0, 1]$ .
- Given the number of data points, the precision of quantile binning can be suboptimal when determining the appropriate bin count.
- By contrast, with local regression, one can specify the percentage of nearest neighbors, providing greater flexibility.

# Our new metric: LCS

## Local Calibration Score (LCS)

A local regression of degree 0, denoted as  $\hat{g}$ , is fitted to the predicted scores  $\hat{s}(\dagger)$ . This fit is then applied to a vector of ~~YCs~~  $l_j$  within the interval  $[0, 1]$ . Each of these points is denoted by  $l_j$ , where  $j \in \{1, \dots, J\}$ , with  $J$  being the target number of points on the visualization curve.

The LCS is defined as:

$$\text{LCS} = \sum_{j=1}^J w_j (\hat{g}(l_j) - l_j)^2, \quad (6)$$

where  $w_j$  is a weight defined as the density of the *score* at  $l_j$ .

# Our new metric: LCS

## Local Calibration Score (LCS)

A local regression of degree 0, denoted as  $\hat{g}$ , is fitted to the predicted scores  $\hat{s}(\dagger)$ . This fit is then applied to a vector of ~~YCs~~  $l_j$  within the interval  $[0, 1]$ . Each of these points is denoted by  $l_j$ , where  $j \in \{1, \dots, J\}$ , with  $J$  being the target number of points on the visualization curve.

The LCS is defined as:

$$\text{LCS} = \sum_{j=1}^J w_j (\hat{g}(l_j) - l_j)^2, \quad (6)$$

where  $w_j$  is a weight defined as the density of the *score* at  $l_j$ .

Note: Austin and Steyerberg, 2019 defined a similar metric using a L1 norm.

# Roadmap

## Impact of Poor Calibration



# Data Generating Process

We  $\mathcal{S} \sim \mathcal{Y} \times \mathcal{C}$  binary observations as in Gutman et al., 2022:

$$D_i \sim \mathcal{B}(p_i),$$

where individual probabilities are obtained using a logistic sigmoid function:

$$p_i = \frac{1}{1 + \exp(-\mathbf{a}^\top \mathbf{x}_i)},$$

$$\eta_i = -\mathbf{a}^\top \mathbf{x}_i + \varepsilon_i$$

with  $\mathbf{a} = [a_1 \ a_2 \ a_3 \ a_4] = [0.1 \ 0.05 \ 0.2 \ -0.05]$  and

$$\mathbf{x}_i = [x_{1,i} \ x_{2,i} \ x_{3,i} \ x_{4,i}]^\top.$$

The observations  $\mathbf{x}_i$  are drawn from a  $\mathcal{U}(0, 1)$  and  $\varepsilon_i \sim \mathcal{N}(0, 0.5^2)$ .

## Forcing Poor Calibration

To simulate  $\hat{y} \leftarrow \mathcal{B}(\mathbf{q}, \mathcal{S})$ , we generate samples of 2,000 observations and we apply (monotonous) transformations to the true probabilities, either on:

- 1 the latent probability  $p_i$ :

$$p_i^u = \left( \frac{1}{1 + 2 \tanh(\mathbb{T} - \eta_i)} \right)^\alpha . \quad (7)$$

- 2 the linear predictor  $\eta_i$ :

$$\eta_i^u = \gamma \times ((-0.1)x_1 + 0.05x_2 + 0.2x_3 - 0.05x_4 + \varepsilon_i) . \quad (8)$$

# Forcing Poor Calibration

To simulate  $\hat{y} \leftarrow \mathcal{B}(\eta; \mathcal{S})$ , we generate samples of 2,000 observations and we apply (monotonous) transformations to the true probabilities, either on:

- 1 the latent probability  $p_i$ :

$$p_i^u = \left( \frac{1}{1 + 2 \tanh(\eta_i)} \right)^\alpha . \quad (7)$$

- 2 the linear predictor  $\eta_i$ :

$$\eta_i^u = \gamma \times ((-0.1)x_1 + 0.05x_2 + 0.2x_3 - 0.05x_4 + \varepsilon_i) . \quad (8)$$

The resulting transformed probabilities are considered as the scores:  $\hat{s}(\dagger) := p_i^u$

## Distortions

- We examine variations in  $\{1/3, 1, 3\}$  for  $\alpha$  and  $\gamma$
- For each of the 6 scenarios, we generate 200 samples of 2,000 obs.

# Distortions

- We examine variations in  $\{1/3, 1, 3\}$  for  $\alpha$  and  $\gamma$
- For each of the 6 scenarios, we generate 200 samples of 2,000 obs.



**Figure 2:** Distorted Prob. as a Function of True Prob., Depending on the Value of  $\alpha$  (left) or  $\gamma$  (right)

# Impacts: Calibration Metrics

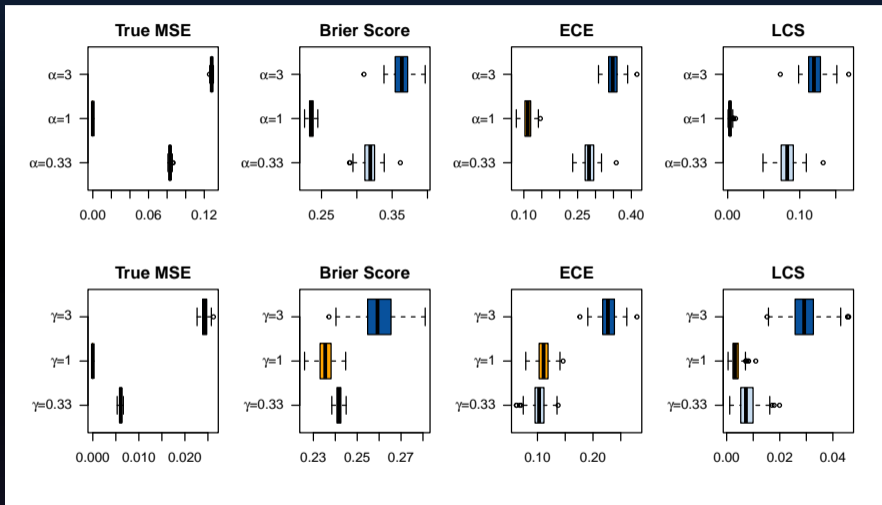
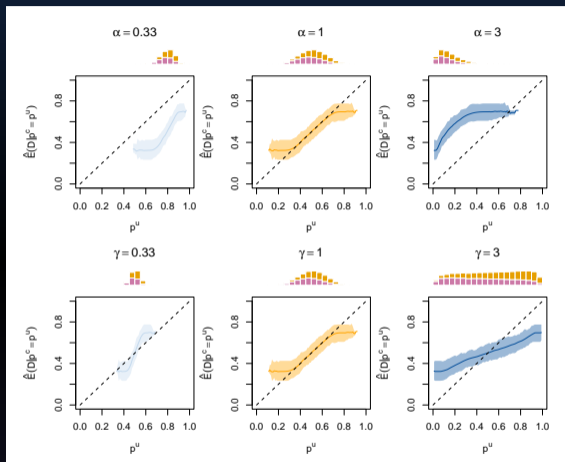


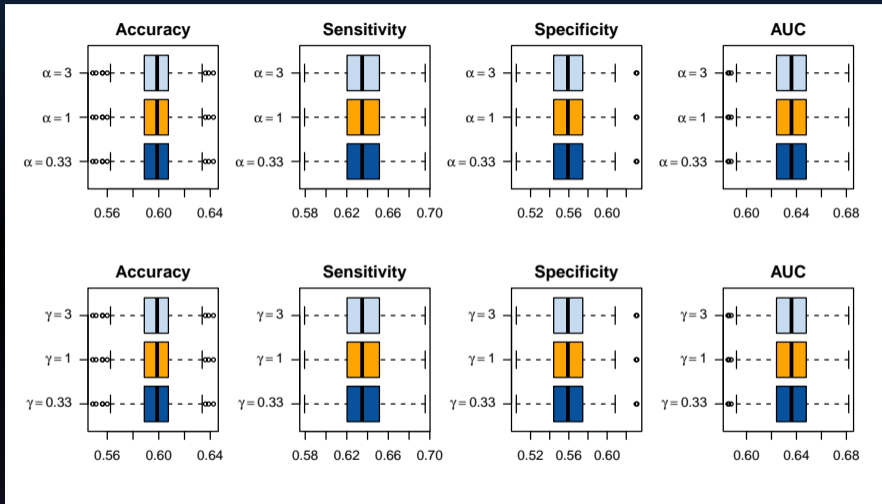
Figure 3: Calibration Metrics on 200 Simulations for each Value of  $\alpha$  (top) or  $\gamma$  (bottom).

# Impacts: Calibration Curves



**Figure 4:** Calibration Curve Obtained with Local Regression, on 200 simulations for each Value of  $\alpha$  (top) or  $\gamma$  (bottom). Distribution of the true probabilities are shown in the histograms (gold for  $d = 1$ , purple for  $d = 0$ ).

# (Mis-)Calibration and standard metrics



**Figure 5:** Standard Goodness of Fit Metrics on 200 Simulations for each Value of  $\alpha$  (top) or  $\gamma$  (bottom). The probability threshold is set to  $\tau = 0.5$ .



# Roadmap

## Calibration and Tree-Based Methods

# Calibration for Machine Learning Algorithms

- With the promise of better performance from machine learning models, it can be tempting to rely on these types of models, such as random forests or derivatives, to estimate binary events.

# Calibration for Machine Learning Algorithms

- With the promise of better performance from machine learning models, it can be tempting to rely on these types of models, such as random forests or derivatives, to estimate binary events.
- However, the distribution of scores returned by these models can be far from the distribution of the underlying probabilities.

# Calibration for Machine Learning Algorithms

- With the promise of better performance from machine learning models, it can be tempting to rely on these types of models, such as random forests or derivatives, to estimate binary events.
- However, the distribution of scores returned by these models can be far from the distribution of the underlying probabilities.
- Here we present an overview of the preliminary results we have obtained with regression trees.

# Calibration for Machine Learning Algorithms

- With the promise of better performance from machine learning models, it can be tempting to rely on these types of models, such as random forests or derivatives, to estimate binary events.
- However, the distribution of scores returned by these models can be far from the distribution of the underlying probabilities.
- Here we present an overview of the preliminary results we have obtained with regression trees.
- More results in the next version of the paper...

# Trees

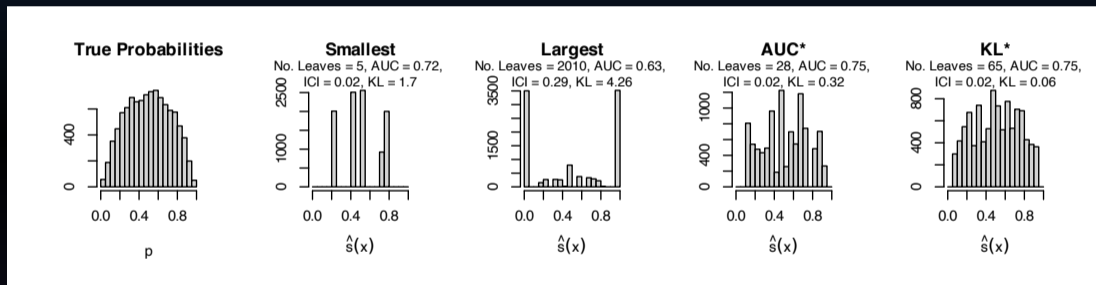
Are trees well calibrated?

# Trees

Are trees well calibrated?

- Some learning algorithms are designed to yield well-calibrated probabilities. These include  $\text{GSS}^{\wedge}$   $\text{zFCs}$ , whose leaf probabilities are optimal on the training set (Kull et al., 2017)
- Earlier studies show that also classical methods such as  $\text{GSS}^{\wedge}$   $\text{zFCs}$ , boosting, SVMs and naive Bayes classifiers tend to be miscalibrated (Wenger et al., 2020)

# Preliminary Results: Calibration is not enough



**Figure 6:** Distribution of true probabilities and estimated scores on validation set for trees of interest,  $n = 10,000$



# Wrap up

- Our new metric, the  $\int_{\mathcal{Y}} |y - \hat{y}| q(\hat{y}) d\mathbb{P}$  offers a more flexible way to visualise and measure calibration than methods based on empirical quantiles.
- $\int_{\mathcal{Y}} |y - \hat{y}| q(\hat{y}) d\mathbb{P}$  : when training classifiers, looking at calibration of models should not be disregarded.

Comments are welcome: [ewen.gallic@univ-amu.fr](mailto:ewen.gallic@univ-amu.fr)

# References I

Austin, P. C. and Steyerberg, E. W. (2019).

The integrated calibration index (ici) and related metrics for quantifying the calibration of logistic regression models.

*Medical Decision Making*, 39(2):140-146, 2019.

Bai, Y., Mei, S., Wang, H., and Xiong, C. (2021).

Don't just blame over-parametrization for over-confidence: Theoretical analysis of calibration in binary classification.

In *Proceedings of the 38th International Conference on Machine Learning*, pages 566-576. PMLR, 2021.

Brier, G. W. (1950).

Verification of forecasts expressed in terms of probability.

*Monthly Weather Review*, 78(1):1-3, 1950.

## References II

Dawid, A. P. (1982).

The well-calibrated bayesian.

*Journal of the Royal Statistical Society B*, 44(2), 187-199.

Gutman, R., Karavani, E., and Shimoni, Y. (2022).

Propensity score models are better when post-calibrated.

Kull, M., Filho, T. M. S., and Flach, P. (2017).

Beyond sigmoids: How to obtain well-calibrated probabilities from binary classifiers with beta calibration.

*Journal of Machine Learning Research*, 18(1), 1-15.

## References III

Loader, C. (1999).

Bayesian Calibration of Probabilities. *Journal of the Royal Statistical Society B*, pages 45–58.

Springer New York, New York, NY.

Pakdaman Naeini, M., Cooper, G., and Hauskrecht, M. (2015).

Obtaining well calibrated probabilities using bayesian binning.

*Journal of Machine Learning Research*, 16:2501–2507.

Schervish, M. J. (1989).

A General Method for Comparing Probability Assessors.

*Journal of the Royal Statistical Society B*, 51(4):1055–1079.

## References IV

Van Calster, B., McLernon, D. J., van Smeden, M., Wynants, L., and Steyerberg, E. W. (2019).

Calibration: the achilles heel of predictive analytics.

"  $J^* = J_2 / B + \beta_{17}(2)$ .

Wenger, J., Kjellström, H., and Triebel, R. (2020).

Non-parametric calibration for classification.

In  $S^Q + 22 / BM; b Q 7 i ? 2 k j` / AM i 2` M i B Q M H * Q M 7 2` 2 M + 2 Q M ` i B } + \beta H A M i$   
 Proceedings of Machine Learning Research.

Wilks, D. S. (1990).

On the combination of forecast probabilities for consecutive precipitation periods.

$q 2 i ? 2` M / 6 Q` 2, 5(4) i 6 14 50$ .

## (Mis-)Calibration and standard metrics

What are the impacts of miscalibration on standard metrics?

We will consider metrics based on the predictive performances calculated using a confusion table:

Table 1: Confusion Table

	Actual Positive	Actual Negative
Predicted Positive	TP	FP
Predicted Negative	FN	TN

where

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}; \quad \text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

# (Mis-)Calibration and standard metrics

$$\text{Accuracy} = \frac{TP + TN}{N}$$

Overall correctness of the model

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

Ability to correctly identify positive class

$$\text{Specificity} = TPR = \frac{TN}{TN + FP}$$

Ability to correctly identify negative class

AUC (Area Under Curve)

TPR and TFP for various prob. threshold  $\tau$