

# Part IV - B

## Model calibration

# Conformal Prediction

Adrian Galdran  
RyC/ikerbasque Research Fellow  
Tecnalia - Derio, Spain  
[adrian.galdran@tecnalia.com](mailto:adrian.galdran@tecnalia.com)



**tecnalia**

MEMBER OF BASQUE RESEARCH  
& TECHNOLOGY ALLIANCE

# Conformal Prediction - Contents

1. Motivation

2. Conformal Prediction: Ingredients

3. Conformal Predictions: Algorithm

4. Hands-On

# 1. Motivation

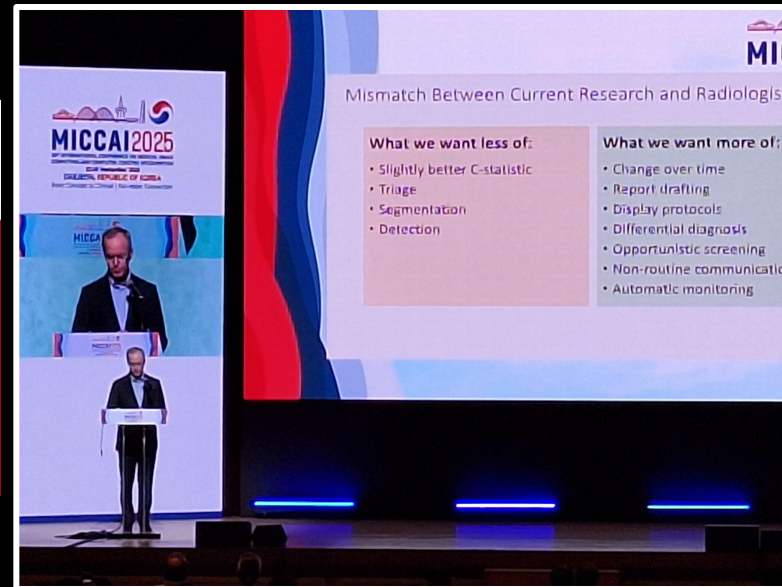


CAP Profiles

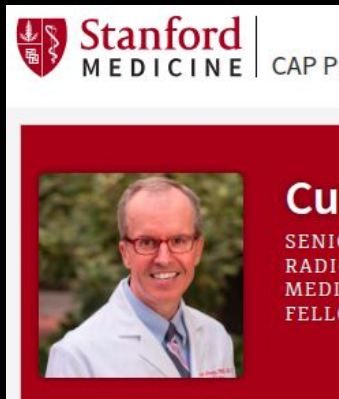


**Curtis Langlotz**

SENIOR ASSOCIATE VICE PROVOST FOR RESEARCH, PROFESSOR OF RADIOLOGY (INTEGRATIVE BIOMEDICAL IMAGING INFORMATICS), OF MEDICINE (BMIR), OF BIOMEDICAL DATA SCIENCE AND SENIOR FELLOW AT THE STANFORD INSTITUTE FOR HUMAN-CENTERED AI



# 1. Motivation



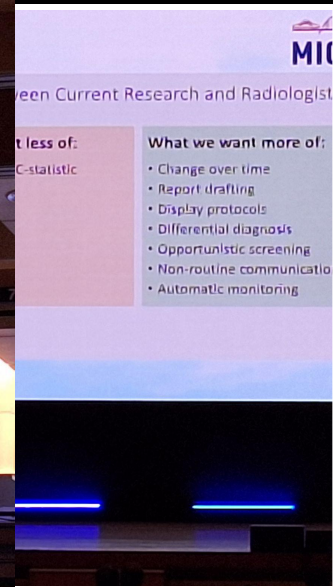
## Mismatch Between Current Research and Radiologist Needs

### What we want less of:

- Slightly better C-statistic
- Triage
- Segmentation
- Detection

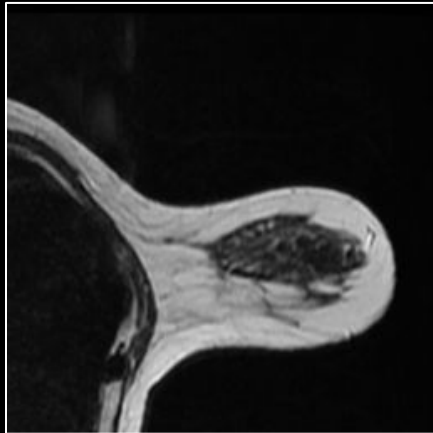
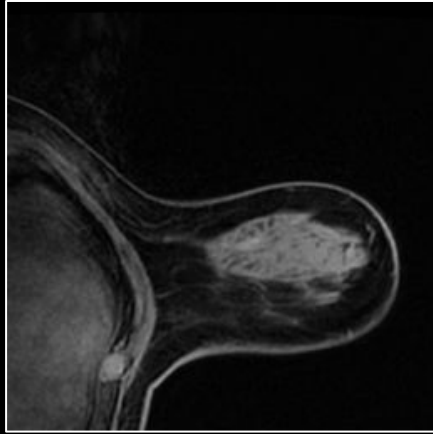
### What we want more of:

- Change over time
- Report drafting
- Display protocols
- **Differential diagnosis**
- Opportunistic screening
- Non-routine communication
- Automatic monitoring

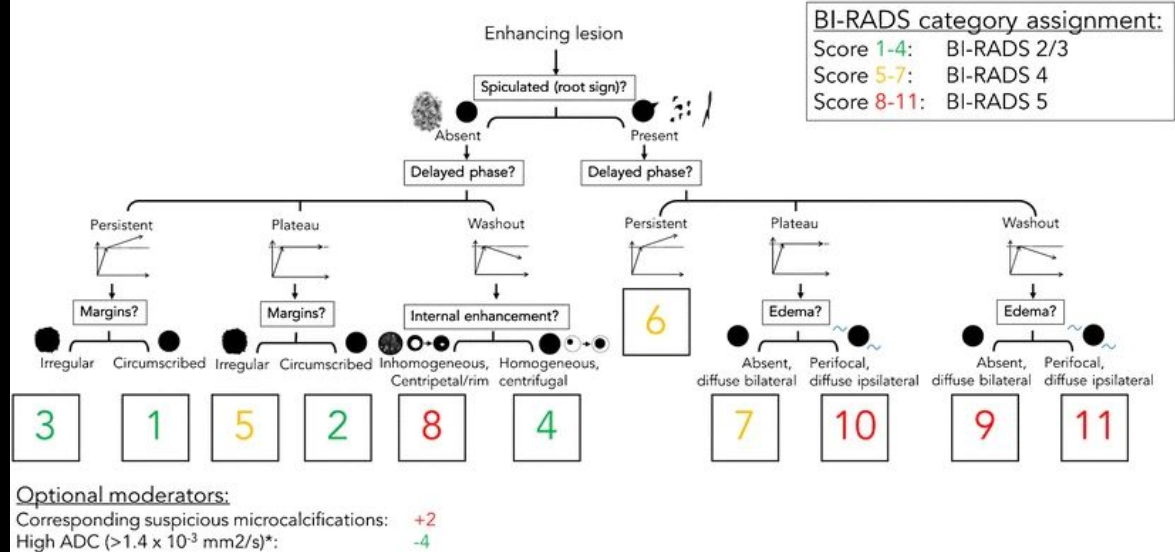


# 1. Motivation

*Dietzel M, Baltzer PAT. How to use the Kaiser score as a clinical decision rule for diagnosis in multiparametric breast MRI: a pictorial essay. Insights Imaging. 2018*

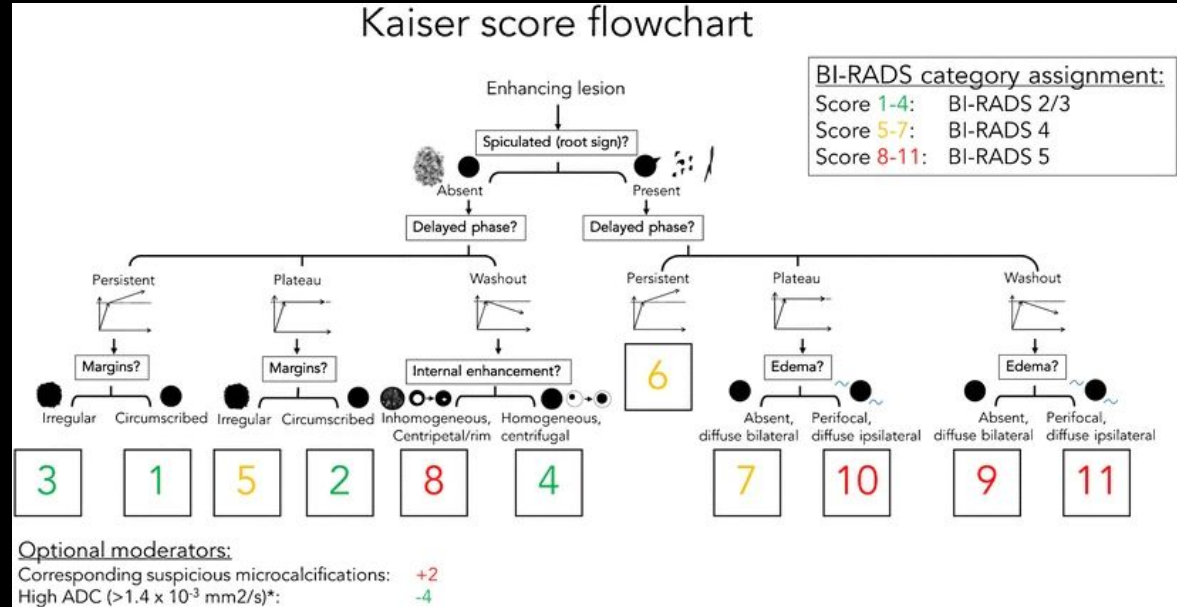


Kaiser score flowchart



# 1. Motivation

*Dietzel M, Baltzer PAT. How to use the Kaiser score as a clinical decision rule for diagnosis in multiparametric breast MRI: a pictorial essay. Insights Imaging. 2018*



$$\mathbb{P}(\text{BI-RADS} \in \{7, 9, 10, 11\}) \geq 90\%$$



# 1. Motivation: Differential Diagnosis

## Step 1: Do we see an enhancing lesion?

- On MRI, when you inject contrast, suspicious lesions often “light up” because they have abnormal blood vessels.
- So the first question is: *does the lesion actually enhance?* If yes → move on.

## Step 2: Is there a spiculated margin (“root sign”)?

- **Spiculation** means the lesion has spikes or radiating lines extending into the surrounding tissue — like roots of a tree.
- This is a strong red flag for cancer. If it’s present, you follow the *right-hand branch*.
- If absent (the lesion looks smooth/rounded), you go to the *left-hand branch*.

## Step 3: What happens in the delayed phase (contrast wash-out curve)?

This is the part you asked about — it’s the **time course** of contrast enhancement.

- After contrast injection, radiologists watch how bright the lesion gets over time.
- There are three typical “curves”:
  - **Persistent:** keeps getting brighter and brighter with time.  
→ Usually benign (think of a sponge slowly soaking water).
  - **Plateau:** gets bright quickly, then levels off.  
→ Suspicious (like tissue that soaks fast but then “caps out”).
  - **Washout:** gets bright early, but then *fades* as contrast drains away.  
→ Very suspicious for cancer (because malignant tumors often have “leaky” vessels).

So the “**delayed phase**” check is: which of these three time-curves does the lesion follow?

## Step 4: If persistent or plateau → check margins

- **Margins** = edge of the lesion.
- Smooth (circumscribed) edges → usually benign.
- Irregular/jagged edges → more worrisome.

## Step 5: If washout → check internal enhancement pattern

- Inside the lesion, how does the contrast distribute?
- If it’s patchy, rim-shaped, or irregular → higher suspicion.
- If it’s uniform or “centrifugal” (from inside out), less worrisome.

## Step 6: If spiculation was present (right side of the tree) → check delayed phase again, then edema

- If the lesion is spiculated and the curve is persistent/plateau/washout, you still branch down.
- In later branches, radiologists look for **edema** — swelling in the tissue around the lesion.
- Edema often shows up in malignancy, so its presence increases the suspicion.

## Putting it together

At the bottom of the tree, you land on a **Kaiser score number (1–11)**.

- **Low scores (1–4):** likely benign (BI-RADS 2/3).
- **Middle scores (5–7):** indeterminate but suspicious (BI-RADS 4).
- **High scores (8–11):** very suspicious, likely malignant (BI-RADS 5).

## 2. Conformal Prediction: Ingredients

Suppose you've got an FDA-approved diagnostic model

$$\hat{\mathcal{M}}_y(x) \sim \mathbb{P}(Y = y | X = x) \quad y \in \{1, \dots, K\} = \mathcal{Y}$$



## 2. Conformal Prediction: Ingredients

Suppose you've got an FDA-approved diagnostic model

$$\hat{\mathcal{M}}_y(x) \sim \mathbb{P}(Y = y | X = x) \quad y \in \{1, \dots, K\} = \mathcal{Y}$$

$$\hat{\mathcal{M}}_y(x^*) = 7$$



## 2. Conformal Prediction: Ingredients

Suppose you've got an FDA-approved diagnostic model

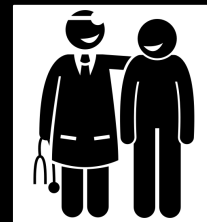
$$\hat{\mathcal{M}}_y(x) \sim \mathbb{P}(Y = y | X = x) \quad y \in \{1, \dots, K\} = \mathcal{Y}$$

$$\hat{\mathcal{M}}_y(x^*) = 7$$



$$\hat{\mathcal{M}}_y(x^*) \in \{7, 8, 9\}$$

*with  $p = 0.9$*



## 2. Conformal Prediction: Ingredients

Suppose you've got an FDA-approved diagnostic model

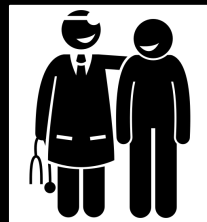
$$\hat{\mathcal{M}}_y(x) \sim \mathbb{P}(Y = y | X = x) \quad y \in \{1, \dots, K\} = \mathcal{Y}$$

$$\hat{\mathcal{M}}_y(x^*) = 7$$



$$\hat{\mathcal{M}}_y(x^*) \in \{7, 8, 9\}$$

*with  $p = 0.9$*



Predict a set  $\mathcal{T}_{x^*} \subseteq \mathcal{Y}$  which contains  $y^*$  *with high  $p$*

$$\text{Coverage : } \mathbb{P}(y^* \in \mathcal{T}_{x^*}) \geq 1 - \alpha$$

## 2. Conformal Prediction: Ingredients

Prediction Sets :  $x^* \mapsto \mathcal{T}_{x^*} \subseteq \mathcal{Y} = \{1, \dots, \mathcal{K}\}$

Coverage :  $\mathbb{P}(y^* \in \mathcal{T}_{x^*}) \geq 1 - \alpha$

---

## 2. Conformal Prediction: Ingredients

Prediction Sets :  $x^* \mapsto \mathcal{T}_{x^*} \subseteq \mathcal{Y} = \{1, \dots, \mathcal{K}\}$

Coverage :  $\mathbb{P}(y^* \in \mathcal{T}_{x^*}) \geq 1 - \alpha$

---



{ fox  
squirrel  
0.99 }



{ fox squirrel, gray  
0.82 fox, bucket, rain  
0.03 0.02 barrel  
0.02 }



{ marmot, fox  
0.30 0.22 squirrel, mink, weasel, beaver, polecat  
0.18 0.16 0.03 0.01 }

## 2. Conformal Prediction: Ingredients

Prediction Sets :  $x^* \mapsto \mathcal{T}_{x^*} \subseteq \mathcal{Y} = \{1, \dots, \mathcal{K}\}$

Coverage :  $\mathbb{P}(y^* \in \mathcal{T}_{x^*}) \geq 1 - \alpha$       Efficient sets

---



{ fox  
squirrel  
0.99 }



{ fox squirrel, gray  
0.82 0.03 bucket, rain  
0.02 barrel  
0.02 }



{ marmot, fox  
0.30 0.22 squirrel, mink, weasel, beaver, polecat  
0.18 0.16 0.03 0.01 }



### 3. Conformal Predictions: Algorithm

Let's get back to our FDA-approved diagnostic model  $\hat{\mathcal{M}}$

### 3. Conformal Predictions: Algorithm

Let's get back to our **FDA-approved** diagnostic model  $\hat{\mathcal{M}}$

**No retraining allowed**

### 3. Conformal Predictions: Algorithm

Let's get back to our **FDA-approved** diagnostic model  $\hat{\mathcal{M}}$

**No retraining allowed**, but you have some fresh data available:

$$(\text{bMRI}_1, \mathbf{y}_1), \dots, (\text{bMRI}_N, \mathbf{y}_i) = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \sim \mathbb{P} \text{ i.i.d.}$$

We call this Calibration Set (sorry about that)

### 3. Conformal Predictions: Algorithm

Let's get back to our **FDA-approved** diagnostic model  $\hat{\mathcal{M}}$

**No retraining allowed**, but you have some fresh data available:

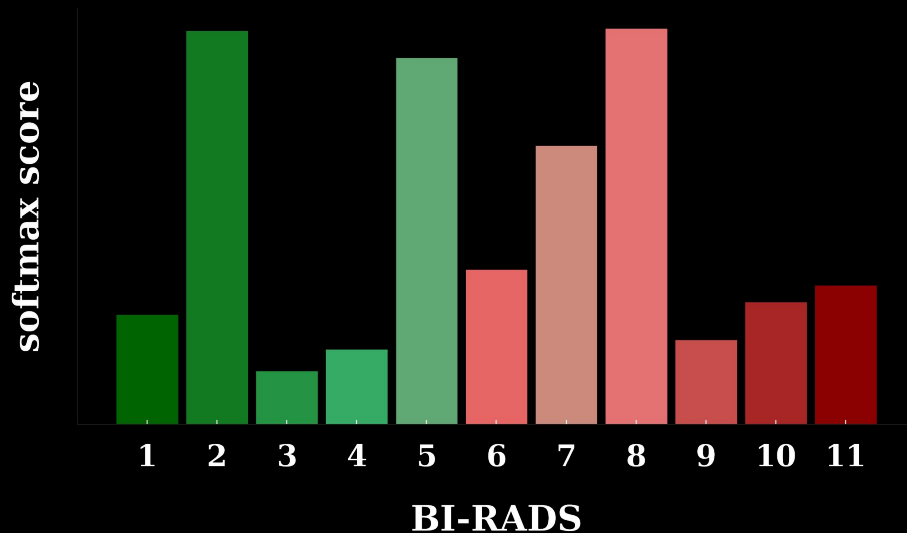
$$(\text{bMRI}_1, \mathbf{y}_1), \dots, (\text{bMRI}_N, \mathbf{y}_i) = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N \sim \mathbb{P} \text{ i.i.d.}$$

We call this Calibration Set (sorry about that)

How do we build these Conformal Prediction Sets?

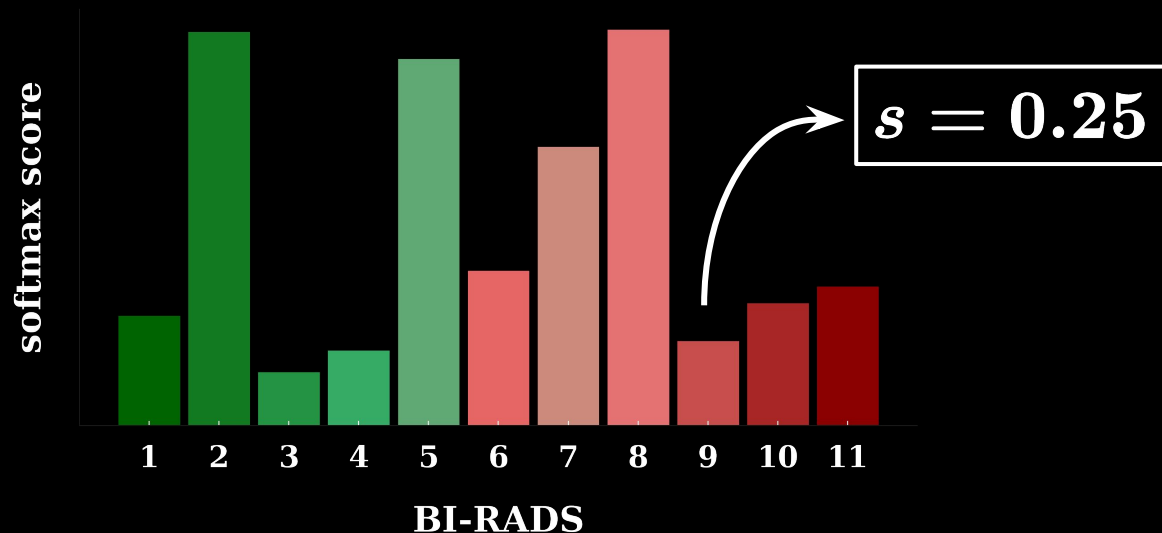
# 3. Conformal Predictions: Algorithm

## 1. Collect scores of correct classes



# 3. Conformal Predictions: Algorithm

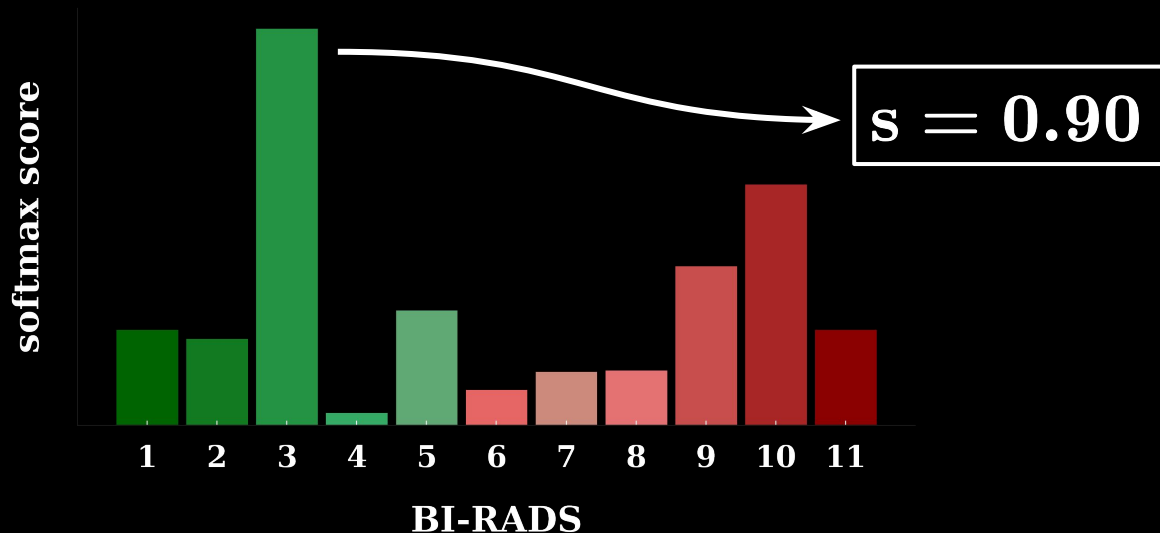
1. Collect scores of correct classes  $E=\{0.25, \}$





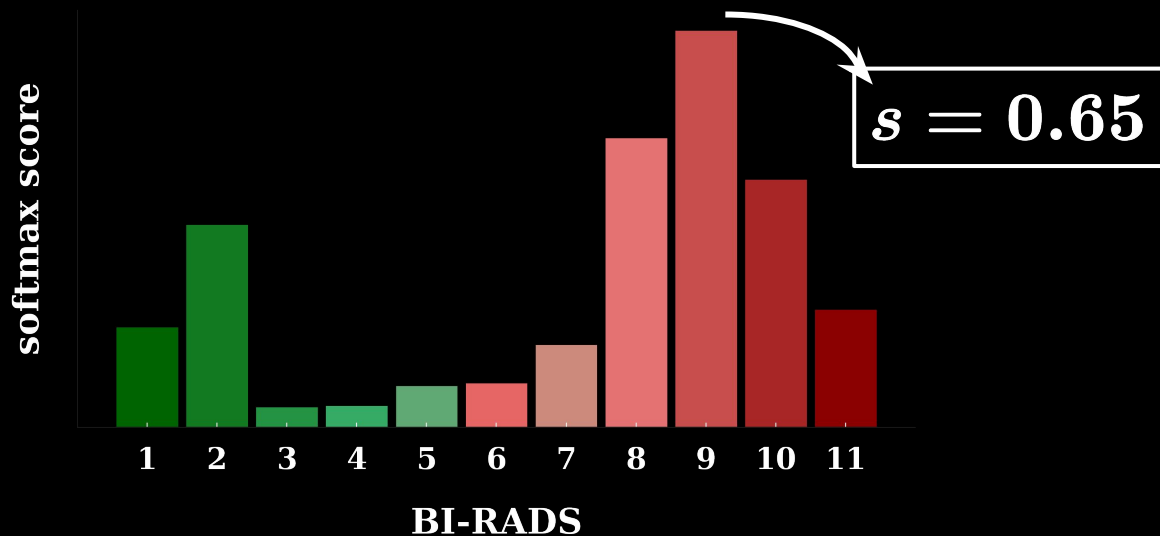
### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $E=\{0.25, 0.90, \}$



### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $E = \{0.25, 0.90, \dots, 0.65\}$



### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $\mathbb{E} = \{0.25, 0.90, \dots, 0.65\}$
2. For a desired coverage of  $\alpha$ , find a value  $\hat{q}_\alpha$  such that you keep  $1-\alpha$  of the scores in  $\mathbb{E}$ :

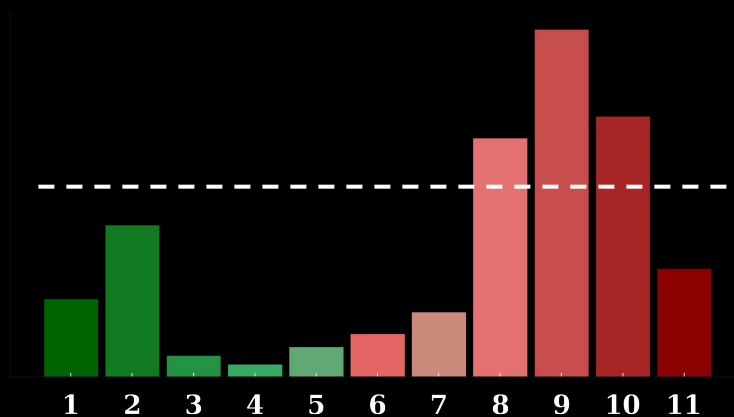
$$\hat{q}_\alpha = \text{np.quantile}([E_1, E_2, \dots, E_N], \alpha)$$

### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $\mathbb{E} = \{0.25, 0.90, \dots, 0.65\}$
2. For a desired coverage of  $\alpha$ , find a value  $\hat{q}_\alpha$  such that you keep  $1-\alpha$  of the scores in  $\mathbb{E}$ :

$$\hat{q}_\alpha = \text{np.quantile}([E_1, E_2, \dots, E_N], \alpha)$$

What happens if we use  $\hat{q}_\alpha$  to build prediction sets in the calibration dataset?



### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $\mathbb{E}=\{0.25, 0.90, \dots, 0.65\}$
2. For a desired coverage of  $\alpha$ , find a value  $\hat{q}_\alpha$  such that you keep  $1-\alpha$  of the scores in  $\mathbb{E}$ :

$$\hat{q}_\alpha = \text{np.quantile}([E_1, E_2, \dots, E_N], \alpha)$$

3. If we use  $\hat{q}_\alpha$  to build prediction sets on test data, we have theoretically guaranteed coverage

### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $\mathbb{E} = \{0.25, 0.90, \dots, 0.65\}$
2. For a desired coverage of  $\alpha$ , find a value  $\hat{q}_\alpha$  such that you keep  $1-\alpha$  of the scores in  $\mathbb{E}$ :

$$\hat{q}_\alpha = \text{np.quantile}([E_1, E_2, \dots, E_N], \alpha)$$

3. If we use  $\hat{q}_\alpha$  to build prediction sets on test data, we have theoretically guaranteed coverage, if data is interchangeable.

$$1-\alpha \leq \mathbb{P}(\mathbf{y}^* \in \mathcal{T}_{x^*}^{\hat{q}_\alpha}) \leq (1-\alpha) + \frac{1}{N+1}$$



### 3. Conformal Predictions: Algorithm

1. Collect scores of correct classes  $\mathbb{E} = \{0.25, 0.90, \dots, 0.65\}$
2. For a desired coverage of  $\alpha$ , find a value  $\hat{q}_\alpha$  such that you keep  $1-\alpha$  of the scores in  $\mathbb{E}$ :

$$\hat{q}_\alpha = \text{np.quantile}([E_1, E_2, \dots, E_N], \alpha)$$

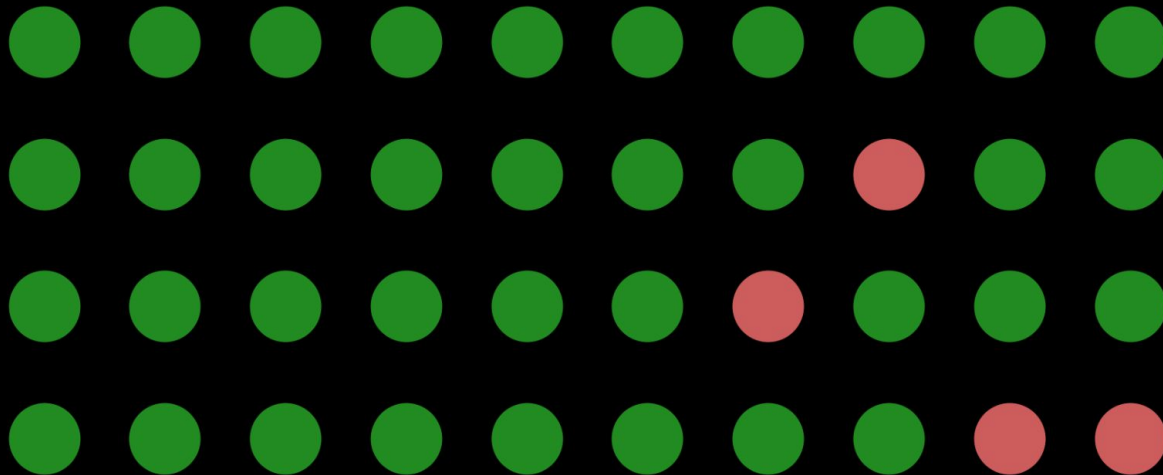
3. If we use  $\hat{q}_\alpha$  to build prediction sets on test data, we have theoretically guaranteed coverage. If data is interchangeable:



$$1-\alpha \leq \mathbb{P}(\mathbf{y}^* \in \mathcal{T}_{x^*}^{\hat{q}_\alpha}) \leq (1-\alpha) + \frac{1}{N+1}$$

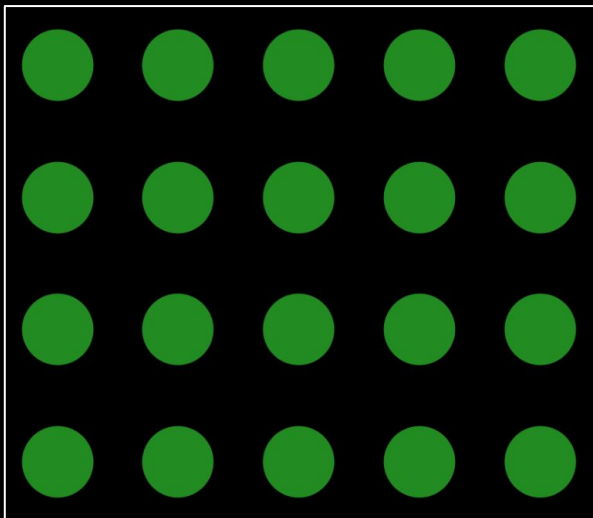
# 3.5 Beyond Coverage

Marginal Coverage of 90%

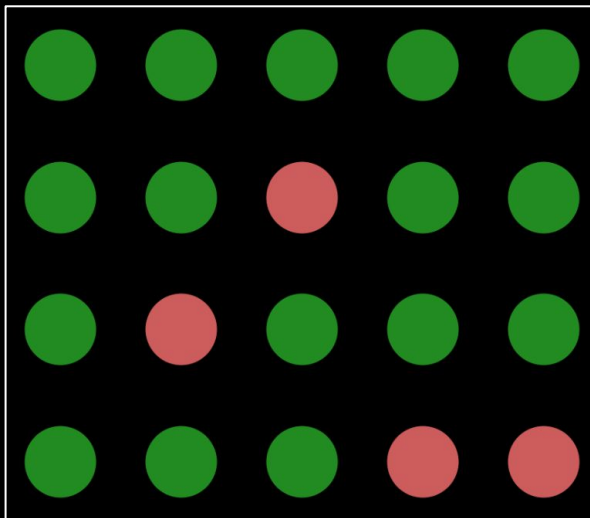


# 3.5 Beyond Coverage

Marginal Coverage of 90%



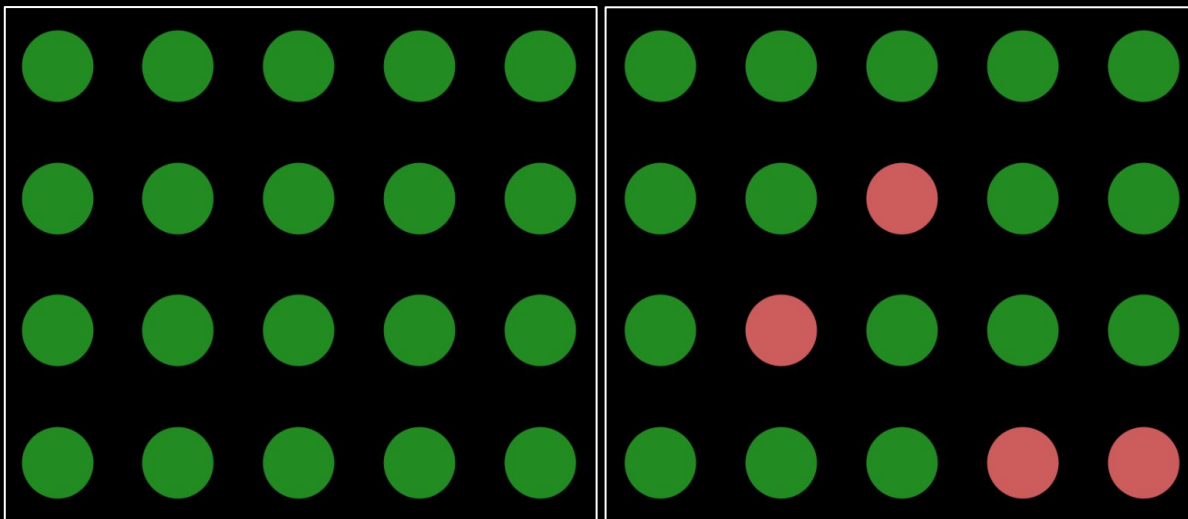
No Conditional Coverage



# 3.5 Beyond Coverage

Marginal Coverage of 90%

No Conditional Coverage



Check out also Conformal Risk Control

## 4. Hands-On

Github repository:

<https://github.com/agaldran/uqinmia-miccai>

