

# Evaluation of fault detection methods in condition monitoring

MOIRA Public Technical Course: 2024/06/07

Douw Marx, Konstantinos Gryllias

# Outline

- 1 Background: Fault detection
- 2 The confusion matrix
- 3 Threshold dependent metric: F1 score
- 4 Threshold independent metrics: ROC and PRC
- 5 The ROC curve and the cost of false alarms
- 6 Conclusions

# Outline

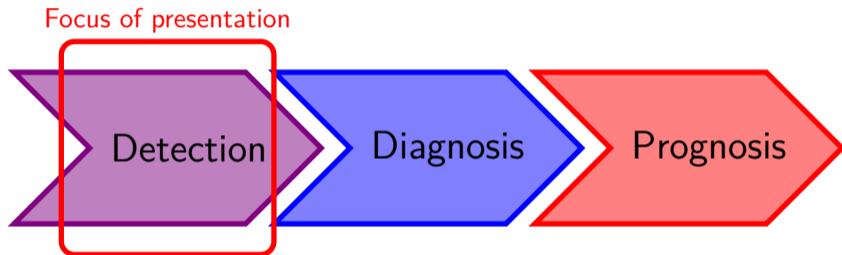
- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms
- ⑥ Conclusions

# Placement of detection in condition monitoring



**Three stages of the condition monitoring**

# Placement of detection in condition monitoring



- ▶ Fault detection
- ▶ Anomaly detection

- ▶ Out of distribution detection
- ▶ Novelty detection

# Motivation for proper fault detection evaluation



**Evaluate the following model** with test data on the left:

**Model:**  $f(\text{apple}) = \text{good}$

(Model that always says an apple is good)

**Accuracy:**

accuracy =  $25/26 = 96\%$  : High accuracy!

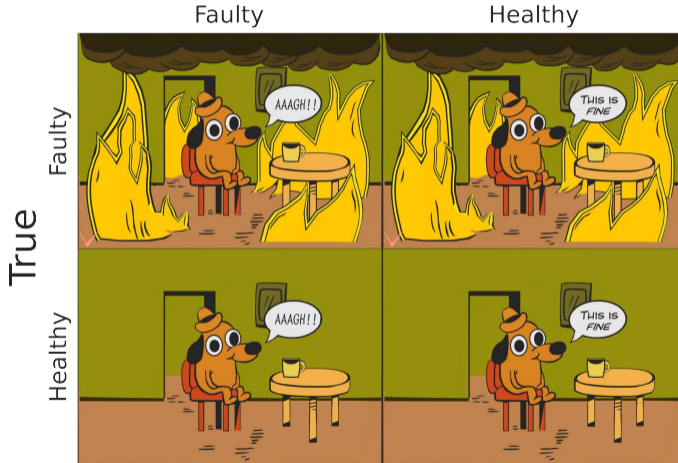
**Is this an effective measure if the model is performing well?**

# Outline

- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms
- ⑥ Conclusions

# What makes for a good fault detection method?

Predicted





# The confusion matrix



Desired Behaviour

# The confusion matrix



Two types of mistakes can be made

# The confusion matrix



Each quadrant has an associated cost

# The confusion matrix

		Predicted	
		Faulty	Not Faulty
Actual	Faulty	TP	FP
	Not Faulty	FN	TN

TP: Correct detection, FP: False alarm,  
FN: Missed detection, TN: Correctly identified as not faulty.

- ▶ Fault detection metrics should incorporate performance on both normal and faulty data (F1, ROC, PRC etc)

# Confusion matrix example

**Model:**  $f(\text{apple}) = \text{good}$  (Model that always says an apple is good)



		Actual	
		Positive	Negative
Predicted	Positive	True Positive 0	False Positive 0
	Negative	False Negative 1	True Negative 25

# Earlier detection is not always better

		Predicted	
		Faulty	Not Faulty
Actual	Faulty	TP	FP
	Not Faulty	FN	TN

TP: Correctly predicted as faulty, FP: Incorrectly predicted as faulty,  
FN: Incorrectly predicted as not faulty, TN: Correctly predicted as not faulty.

- ▶ Early detection is often mainly concerned with **correct detections**.

# Earlier detection is not always better

		Predicted	
		Faulty	Not Faulty
Actual	Faulty	TP	FP
	Not Faulty	FN	TN

TP: Correctly predicted as faulty, FP: Incorrectly predicted as faulty,  
FN: Incorrectly predicted as not faulty, TN: Correctly predicted as not faulty.

- ▶ Early detection is often mainly concerned with **correct detections**.
- ▶ The importance of identifying **reference** samples is often overlooked.

## Earlier detection is not always better

		Predicted	
		Faulty	Not Faulty
Actual	Faulty	TP	FP
	Not Faulty	FN	TN

TP: Correctly predicted as faulty, FP: Incorrectly predicted as faulty,  
FN: Incorrectly predicted as not faulty, TN: Correctly predicted as not faulty.

- ▶ Early detection is often mainly concerned with **correct detections**.
- ▶ The importance of identifying **reference** samples is often overlooked.
- ▶ total cost =  $\text{cost}(\text{TP}) \cdot \text{TP} + \text{cost}(\text{FN}) \cdot \text{FN} + \text{cost}(\text{FP}) \cdot \text{FP} + \text{cost}(\text{TN}) \cdot \text{TN}$



# Outline

- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score**
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms
- ⑥ Conclusions

# Confusion Matrix Cocktails

“Accuracy on positive class”

[sensitivity](#), [recall](#), [hit rate](#), or [true positive rate](#) (TPR)

$$\text{TPR} = \frac{\text{TP}}{\text{P}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = 1 - \text{FNR}$$

[specificity](#), [selectivity](#) or [true negative rate](#) (TNR)

$$\text{TNR} = \frac{\text{TN}}{\text{N}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

“How often was the model correct when it was betting on the positive class”

[precision](#) or [positive predictive value](#) (PPV)

$$\text{PPV} = \frac{\text{TP}}{\text{TP} + \text{FP}} = 1 - \text{FDR}$$

[negative predictive value](#) (NPV)

$$\text{NPV} = \frac{\text{TN}}{\text{TN} + \text{FN}} = 1 - \text{FOR}$$

miss rate or [false negative rate](#) (FNR)

$$\text{FNR} = \frac{\text{FN}}{\text{P}} = \frac{\text{FN}}{\text{FN} + \text{TP}} = 1 - \text{TPR}$$

[fall-out](#) or [false positive rate](#) (FPR)

$$\text{FPR} = \frac{\text{FP}}{\text{N}} = \frac{\text{FP}}{\text{FP} + \text{TN}} = 1 - \text{TNR}$$

[false discovery rate](#) (FDR)

$$\text{FDR} = \frac{\text{FP}}{\text{FP} + \text{TP}} = 1 - \text{PPV}$$

[false omission rate](#) (FOR)

$$\text{FOR} = \frac{\text{FN}}{\text{FN} + \text{TN}} = 1 - \text{NPV}$$

# F1 score for balanced performance

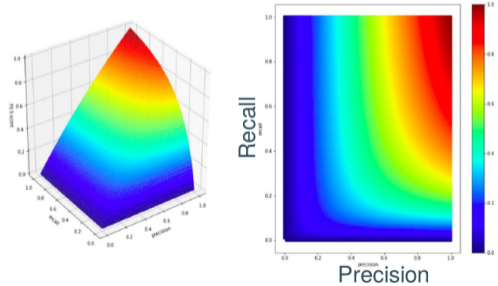
- We want to find a compromise between precision (If model says the data point is positive it is actually positive) and recall/sensitivity (You have a high accuracy on the positive class).

- Measure between 0 and 1

$$F_1 = \frac{2}{\text{recall}^{-1} + \text{precision}^{-1}} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

## Example

- True positives (TP): 75 correctly identified as faulty machinery.
- False positives (FP): 10 samples incorrectly identified as faulty machinery.
- False negatives (FN): 5 faulty machinery samples incorrectly classified as non-faulty.
- True negatives (TN): 10 samples were correctly identified as non-faulty machinery.
  
- Precision = TP / (TP + FP) = 75/(75+10) = 0.882
- Recall = TP / (TP + FN) = 75/(75+5) = 0.938
  
- F1 score = 2 \* ((Precision \* Recall) / (Precision + Recall)) = 2 \* ((0.882 \* 0.938) / (0.882 + 0.938)) = 0.909

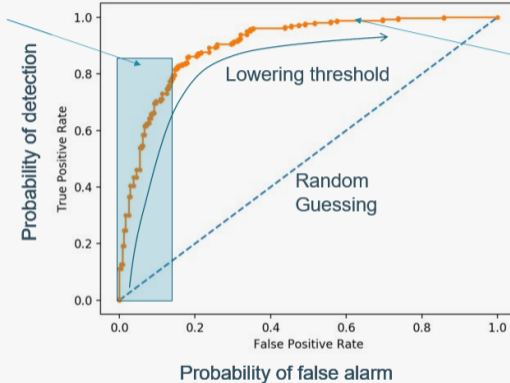


# Outline

- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms
- ⑥ Conclusions

# The Receiver Operating Curve (ROC)

We mostly care about this region for fault detection

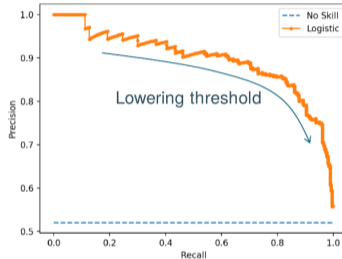


We can compute the area under this curve as summary metric for how well a model is doing over all thresholds.

The ROC curve shows the trade-off between hit rate (TPR) and false alarm rate (FPR).

# Precision Recall Curve (PRC)

- Show the trade-off between precision and recall/sensitivity as the threshold of detection is varied.



Threshold example:

Faulty if RMS vibration above 8G vs 9G vs 10G.

The PRC curve shows the trade-off between precision and recall.

Jason Brownlee, ROC Curves and Precision-Recall Curves for Imbalanced Classification, Machine Learning Mastery, Available from <https://machinelearningmastery.com/roc-curves-and-precision-recall-curves-for-imbalanced-classification/>, accessed May 3rd, 2020.

# Precision Recall Curve vs ROC Curve

Precision Recall Curve:

- ▶ Better for imbalanced datasets: Does not accounts for TN unlike ROC.
- ▶ Less intuitive than ROC

ROC Curve:

- ▶ Often easier to interpret than PRC
- ▶ Overly optimistic for highly imbalanced datasets

# Outline

- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms**
- ⑥ Conclusions



# The ROC curve and CBM cost

$$\text{Cost} = \text{TPR} \times C_{\text{TP}} + \text{TNR} \times C_{\text{TN}} \\ + \text{FPR} \times C_{\text{FP}} + \text{FNR} \times C_{\text{FN}}$$

▶ Assume

- $C_{\text{TP}} = 0$  (No cost for correctly detecting a fault)
- $C_{\text{TN}} = 0$  (No cost for correctly detecting a healthy sample)

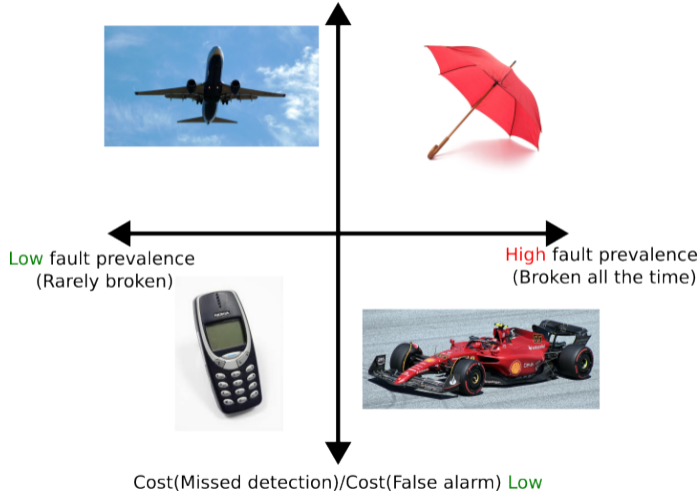
▶  $C_{\text{FP}}$  (False alarm cost) and  $C_{\text{FN}}$  (Missed detection cost) are varied.

▶  $C_X$ : Accounts for 1) Cost per occurrence 2) Prevalence of the fault

▶  $\text{FNR} = 1 - \text{TPR}$

# Cost and prevalence differences for different applications

Cost(Missed detection)/Cost(False alarm) **High**



Both the relative cost of false alarms and missed detections and the prevalence of faults should be considered

# Cost and prevalence differences for different applications

Cost(Missed detection)/Cost(False alarm) **High**



FP: Unnecessary maintenance  
FN: Loss of life



FP: Annoying umbrella notification  
FN: You get soaked

**Low** fault prevalence  
(Rarely broken)



FP: Annoying sms notification  
FN: Can't play snake on bus

**High** fault prevalence  
(Broken all the time)

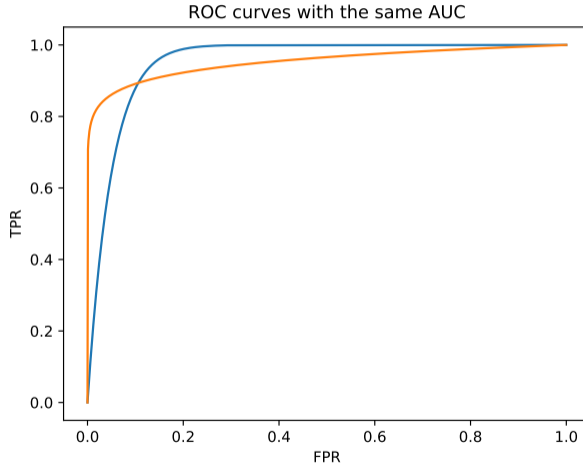


FP: Require unnecessary pit stop  
FN: Complete race with 3 wheels

Cost(Missed detection)/Cost(False alarm) **Low**

Both the relative cost of false alarms and missed detections and the prevalence of faults should be considered

# Models suitable for different applications



Not all ROC curves are created equal.  
Orange: Preferred for e.g. package inspection,  
Blue: Preferred for e.g. nuclear power plants.

# Goal: Optimise ROC curve given cost and prevalence

- ▶ Increase correct detections, without contributing to more false alarms.
- ▶ Consider the relative cost of false alarms and missed detections and the prevalence of faults.

# Outline

- ① Background: Fault detection
- ② The confusion matrix
- ③ Threshold dependent metric: F1 score
- ④ Threshold independent metrics: ROC and PRC
- ⑤ The ROC curve and the cost of false alarms
- ⑥ Conclusions

# Conclusions

- ▶ Evaluation metrics used to evaluate fault detection methods should incorporate performance on both normal and faulty data.
- ▶ Evaluation metrics should be designed based on the application and relative cost of false positives and false negatives.

Thank you for your attention  
Any comments / critique will be appreciated

[douw.marx@kuleuven.be](mailto:douw.marx@kuleuven.be)

The authors gratefully acknowledge the European Commission for its support of the Marie Skłodowska Curie program through the ETN MOIRA project (GA 955681).