

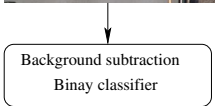
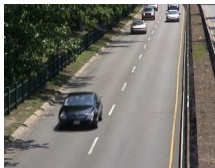
Summarizing the performances of a background subtraction algorithm measured on several videos

Sébastien Piérard and [Marc Van Droogenbroeck](#)

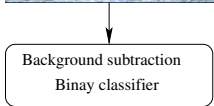
Department of Electrical Engineering and Computer Science (Montefiore Institute), University of Liège, Belgium

Special Session on “Dynamic Background Reconstruction/Subtraction for Challenging Environments”

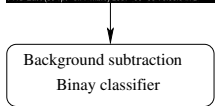
Motivation: scoring an algorithm for multiple videos



$P_1, R_1, TPR_1, ER_1, F_1$

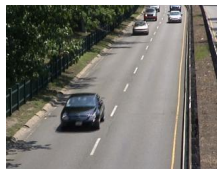


$P_2, R_2, TPR_2, ER_2, F_2$



$P_3, R_3, TPR_3, ER_3, F_3$

Scoring multiple videos with a **unique series of indicators**



Background subtraction
Binay classifier



Background subtraction
Binay classifier



Background subtraction
Binay classifier



P, R, TPR, ER, F

Should we use the mean for scoring multiple videos?

Do we have a candidate mechanism for aggregating scores of multiple videos?

A natural/obvious candidate for scoring multiple videos is the (arithmetic) mean.

So, if we have:

- ▶ Performance for video 1: P_1, R_1, F_1
- ▶ Performance for video 2: P_2, R_2, F_2

we calculate

$$\bar{P} = \frac{P_1 + P_2}{2}, \bar{R} = \frac{R_1 + R_2}{2}, \text{ and } \bar{F} = \frac{F_1 + F_2}{2}$$

But there is a catch ...

Should we use the mean for scoring multiple videos?

Do we have a candidate mechanism for aggregating scores of multiple videos?

A natural/obvious candidate for scoring multiple videos is the (arithmetic) mean.

So, if we have:

- ▶ Performance for video 1: P_1, R_1, F_1
- ▶ Performance for video 2: P_2, R_2, F_2

we calculate

$$\bar{P} = \frac{P_1 + P_2}{2}, \bar{R} = \frac{R_1 + R_2}{2}, \text{ and } \bar{F} = \frac{F_1 + F_2}{2}$$

But there is a catch ...

We should not use the mean for scoring multiple videos!

Obviously,

$$\text{for any video } i, F_i = 2 \frac{P_i \times R_i}{P_i + R_i} \quad \text{but} \quad \bar{F} \neq 2 \frac{\bar{P} \times \bar{R}}{\bar{P} + \bar{R}} = \bar{\bar{F}}$$

We should not use the mean for scoring multiple videos!

Obviously,

$$\text{for any video } i, F_i = 2 \frac{P_i \times R_i}{P_i + R_i} \quad \text{but} \quad \bar{F} \neq 2 \frac{\bar{P} \times \bar{R}}{\bar{P} + \bar{R}} = \bar{\bar{F}}$$

The M4CD algorithm of CDNet 2014 typically illustrates the problem

$$\bar{F} = 0.69 \quad \neq \quad \bar{\bar{F}} = 2 \frac{\bar{P} \times \bar{R}}{\bar{P} + \bar{R}} = 0.75$$

In fact, the arithmetic mean has severe drawbacks:

- ▶ it breaks the intrinsic relationships between probabilistic indicators.
- ▶ because of these inconsistencies, we might have that

$$\bar{F}_1 < \bar{F}_2$$

while

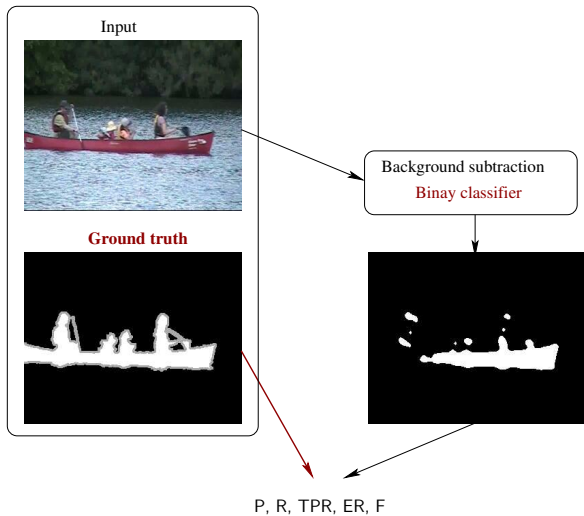
$$\bar{\bar{F}}_1 > \bar{\bar{F}}_2 .$$

Outline

- 1 Performance indicators for one video
- 2 Summarizing the performance for several videos
- 3 Summarizing applied on CDNET 2014
- 4 Conclusion

A scenario for the evaluation of background subtraction algorithms

Dataset



Towards performance indicators applicable to a binary classifier

Ground truth



$y \in \{c^+ = \text{foreground}, c^- = \text{background}\}$

Prediction



$\hat{y} \in \{c^+ = \text{foreground}, c^- = \text{background}\}$

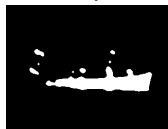
P, R, TPR, ER, F

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
		TP	FN
		FP	TN

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
		TP	FN
		FP	TN

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

The confusion matrix

Ground truth



Output



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

Experimental performance indicators based on the confusion matrix I



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

Positive prior $\pi^+ = \frac{TP+FN}{TP+FN+FP+TN}$

Precision $P = \frac{TP}{TP+FP} = \text{PPV}$ Positive Predictive Value

True Positive Rate $\text{TPR} = \frac{TP}{TP+FN} = R$ Recall

Experimental performance indicators based on the confusion matrix II



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	TP	FN
	Negative	FP	TN

$$\text{Accuracy } A = \frac{TP+TN}{TP+FN+FP+TN}$$

$$\text{Error rate } ER = \frac{FP+FN}{TP+FN+FP+TN}$$

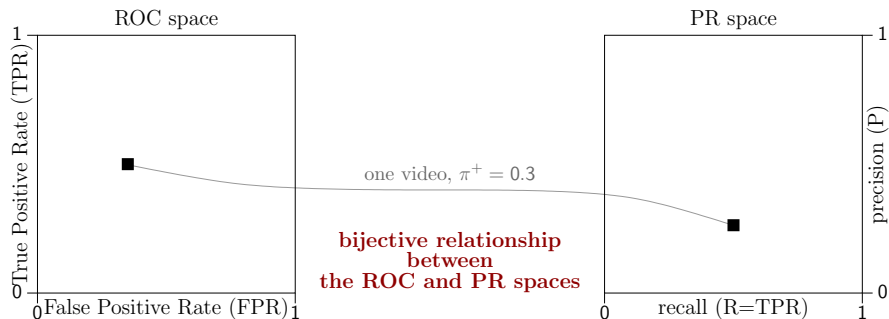
$$\text{F score } F = \frac{2TP}{2TP+FN+FP}$$

ROC vs PR evaluation spaces: there is a bijection!

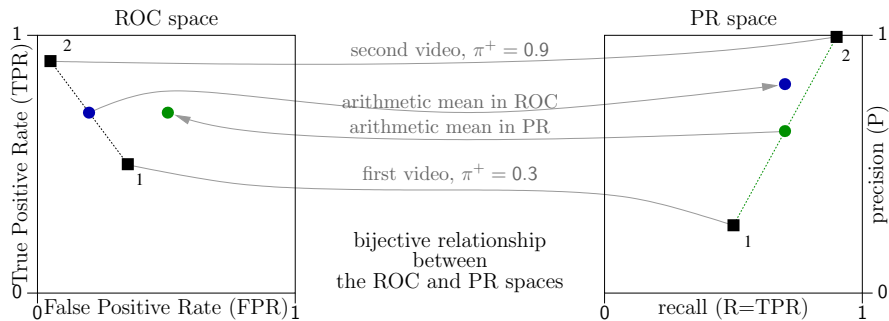
There are two well-known evaluation spaces:

ROC: Receiver Operating Characteristic, defined by (FPR, TPR)

PR: Precision/Recall



Effect of the arithmetic mean



There is no bijection between the means anymore!

The “normalized” confusion matrix



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	pTP	pFN
	Negative	pFP	pTN

The proportion of TP, denoted by pTP, is defined as

$$\frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP} + \text{TN}}$$

This has no impact on the calculation of indicators, such as the F score:

$$F = \frac{2\text{TP}}{2\text{TP} + \text{FN} + \text{FP}} = \frac{2\text{pTP}}{2\text{pTP} + \text{pFN} + \text{pFP}}$$

but it leads to a helpful interpretation of experimental indicators in terms of probabilities.

The “normalized” confusion matrix



		Predicted class \hat{y}	
		Positive	Negative
Actual class y	Positive	pTP	pFN
	Negative	pFP	pTN

The proportion of TP, denoted by pTP, is defined as

$$\frac{TP}{TP + FN + FP + TN}$$

This has no impact on the calculation of indicators, such as the F score:

$$F = \frac{2TP}{2TP + FN + FP} = \frac{2pTP}{2pTP + pFN + pFP}$$

but it leads to a helpful interpretation of experimental indicators in terms of probabilities.

Probabilistic meaning of experimental performance indicators

Definition (Joint random experiment for one video)

Draw one pixel at random (all pixels being equally likely) from the video and jointly observe the ground-truth class Y and the predicted class \hat{Y} for this pixel.

Joint random experiment $\Delta = (Y, \hat{Y})$		Prediction \hat{Y}	
		Positive	Negative
Ground truth Y	Positive	tp = (c^+, c^+)	fn = (c^+, c^-)
	Negative	fp = (c^-, c^+)	tn = (c^-, c^-)

There are four possible outcomes: $\{tp, fn, fp, tn\}$.

Probabilistic indicators

Joint random experiment $\Delta = (Y, \hat{Y})$		Prediction \hat{Y}	
		Positive	Negative
Ground truth Y	Positive	tp = (c^+, c^+)	fn = (c^+, c^-)
	Negative	fp = (c^-, c^+)	tn = (c^-, c^-)

The family of *probabilistic indicators* can be defined based on this random experiment:

$$P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}) \text{ with } \emptyset \subsetneq \mathcal{A} \subsetneq \mathcal{B} \subseteq \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\} \quad (1)$$

It includes

- ▶ $\pi^+ = P(\Delta \in \{\text{tp}, \text{fn}\} | \Delta \in \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\}) = P(\Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $\text{TPR} = R = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $P = \text{PPV} = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fp}\})$, $\text{ER} = P(\Delta \in \{\text{fn}, \text{fp}\})$
- ▶ ... but not the F score!

Probabilistic indicators

Joint random experiment $\Delta = (Y, \hat{Y})$		Prediction \hat{Y}	
		Positive	Negative
Ground truth Y	Positive	tp = (c^+, c^+)	fn = (c^+, c^-)
	Negative	fp = (c^-, c^+)	tn = (c^-, c^-)

The family of *probabilistic indicators* can be defined based on this random experiment:

$$P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}) \text{ with } \emptyset \subsetneq \mathcal{A} \subsetneq \mathcal{B} \subseteq \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\} \quad (1)$$

It includes

- ▶ $\pi^+ = P(\Delta \in \{\text{tp}, \text{fn}\} | \Delta \in \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\}) = P(\Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $\text{TPR} = \text{R} = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $\text{P} = \text{PPV} = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fp}\})$, $\text{ER} = P(\Delta \in \{\text{fn}, \text{fp}\})$
- ▶ ... but not the F score!

Probabilistic indicators

Joint random experiment $\Delta = (Y, \hat{Y})$		Prediction \hat{Y}	
		Positive	Negative
Ground truth Y	Positive	tp = (c^+, c^+)	fn = (c^+, c^-)
	Negative	fp = (c^-, c^+)	tn = (c^-, c^-)

The family of *probabilistic indicators* can be defined based on this random experiment:

$$P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}) \text{ with } \emptyset \subsetneq \mathcal{A} \subsetneq \mathcal{B} \subseteq \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\} \quad (1)$$

It includes

- ▶ $\pi^+ = P(\Delta \in \{\text{tp}, \text{fn}\} | \Delta \in \{\text{tp}, \text{fn}, \text{fp}, \text{tn}\}) = P(\Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $\text{TPR} = \text{R} = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fn}\})$
- ▶ $\text{P} = \text{PPV} = P(\Delta = \text{tp} | \Delta \in \{\text{tp}, \text{fp}\})$, $\text{ER} = P(\Delta \in \{\text{fn}, \text{fp}\})$
- ▶ ... but not the F score!

Outline

- 1 Performance indicators for one video
- 2 Summarizing the performance for several videos**
- 3 Summarizing applied on CDNET 2014
- 4 Conclusion

A probabilistic model for summarization

Definition (Parametric random experiment for several videos)

First, draw one video V at random in the set \mathbb{V} , following an **arbitrarily chosen distribution** $P(V)$. Then, draw one pixel at random from V and observe the ground-truth class Y and the predicted class \hat{Y} for this pixel.

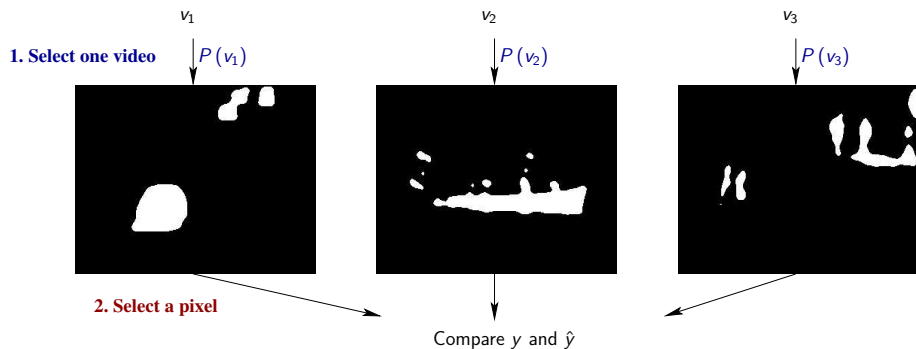


Figure: A probabilistic model for summarization: $\Delta = (V, Y, \hat{Y})$.

Summarization formulas

Notations:

- ▶ $I(v)$ = the value of a performance indicator I for a video $v \in \mathbb{V}$,
- ▶ $I(\mathbb{V})$ = the value of I for a set \mathbb{V} of videos.

We define a probabilistic indicator $I_{\mathcal{A}|\mathcal{B}}$ as $P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B})$, and $I_{\mathcal{B}}$ as $P(\Delta \in \mathcal{B})$. We have

$$\begin{aligned} I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) &= P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}) \\ &= \sum_{v \in \mathbb{V}} P(\Delta \in \mathcal{A}, V = v | \Delta \in \mathcal{B}) \\ &= \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}, V = v) \\ I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) &= \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) I_{\mathcal{A}|\mathcal{B}}(v) \end{aligned} \quad (2)$$

For the particular case of an unconditional probabilistic indicator $I_{\mathcal{A}} = I_{\mathcal{A}|\{\text{tn,fp,fn,tp}\}}$, we have

$$I_{\mathcal{A}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{A}}(v) \quad (3)$$

Summarization formulas

Notations:

- ▶ $I(v)$ = the value of a performance indicator I for a video $v \in \mathbb{V}$,
- ▶ $I(\mathbb{V})$ = the value of I for a set \mathbb{V} of videos.

We define a probabilistic indicator $I_{\mathcal{A}|\mathcal{B}}$ as $P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B})$, and $I_{\mathcal{B}}$ as $P(\Delta \in \mathcal{B})$. We have

$$\begin{aligned} I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) &= P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}) \\ &= \sum_{v \in \mathbb{V}} P(\Delta \in \mathcal{A}, V = v | \Delta \in \mathcal{B}) \\ &= \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) P(\Delta \in \mathcal{A} | \Delta \in \mathcal{B}, V = v) \\ I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) &= \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) I_{\mathcal{A}|\mathcal{B}}(v) \end{aligned} \quad (2)$$

For the particular case of an unconditional probabilistic indicator

$I_{\mathcal{A}} = I_{\mathcal{A}|\{\text{tn,fp,fn,tp}\}}$, we have

$$I_{\mathcal{A}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{A}}(v) \quad (3)$$

Summarization formulas and properties

Formulas:

$$I_{\mathcal{A}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{A}}(v)$$

$$I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) I_{\mathcal{A}|\mathcal{B}}(v)$$

Example: $TPR(\mathbb{V}) = \frac{1}{\pi^+(\mathbb{V})} \sum_{v \in \mathbb{V}} P(V = v) \pi^+(v) TPR(v)$ (4)

Properties:

- 1 Summarization preserves the consistency between indicators, including the bijection between the ROC and PR spaces!
- 2 As long as an indicator is defined for at least one video, it can be summarized! To prove it, we rewrite $I_{\mathcal{A}|\mathcal{B}}(\mathbb{V})$ as

$$I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) = \frac{I_{\mathcal{A} \cap \mathcal{B}}(\mathbb{V})}{I_{\mathcal{B}}(\mathbb{V})} = \frac{I_{\mathcal{A} \cap \mathcal{B}}(\mathbb{V})}{\sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{B}}(v)} \quad (5)$$

Summarization formulas and properties

Formulas:

$$I_{\mathcal{A}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{A}}(v)$$

$$I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) = \sum_{v \in \mathbb{V}} P(V = v | \Delta \in \mathcal{B}) I_{\mathcal{A}|\mathcal{B}}(v)$$

Example: $TPR(\mathbb{V}) = \frac{1}{\pi^+(\mathbb{V})} \sum_{v \in \mathbb{V}} P(V = v) \pi^+(v) TPR(v)$ (4)

Properties:

- 1 Summarization preserves the consistency between indicators, including the bijection between the ROC and PR spaces!
- 2 As long as an indicator is defined for at least one video, it can be summarized! To prove it, we rewrite $I_{\mathcal{A}|\mathcal{B}}(\mathbb{V})$ as

$$I_{\mathcal{A}|\mathcal{B}}(\mathbb{V}) = \frac{I_{\mathcal{A} \cap \mathcal{B}}(\mathbb{V})}{I_{\mathcal{B}}(\mathbb{V})} = \frac{I_{\mathcal{A} \cap \mathcal{B}}(\mathbb{V})}{\sum_{v \in \mathbb{V}} P(V = v) I_{\mathcal{B}}(v)} \quad (5)$$

An algorithm for the computation of summarized indicators

Algorithm:

- 1 Blend the **normalized** confusion matrices with the $P(v_1), P(v_2), \dots$ weights,
- 2 then calculate the indicators!

$$\begin{array}{cc} pTP_1 & pFN_1 & pTP_2 & pFN_2 & pTP_3 & pFN_3 \\ pFP_1 & pTN_1 & pFP_2 & pTN_2 & pFP_3 & pTN_3 \end{array}$$

$$P(v_1)$$

$$P(v_2)$$

$$P(v_3)$$

$$\begin{array}{cc} \sum_v P(v) pTP_v & \sum_v P(v) pFN_v \\ \sum_v P(v) pFP_v & \sum_v P(v) pTN_v \end{array}$$

$$F = \frac{2pTP}{2pTP+pFN+pFP}$$

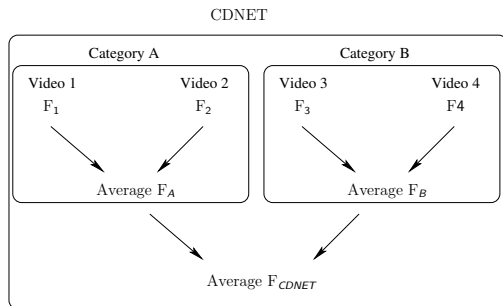
Outline

- 1 Performance indicators for one video
- 2 Summarizing the performance for several videos
- 3 Summarizing applied on CDNET 2014**
- 4 Conclusion

Experiments with CDNET 2014

We analyze two scenarios:

The original CDNET procedure



Our summarization, with

$$P(V = v) = \frac{1}{11} \times \frac{1}{M}$$

$$\begin{array}{ccccc} pTP_1 & pFN_1 & pTP_2 & pFN_2 & pTP_3 & pFN_3 \\ pFP_1 & pTN_1 & pFP_2 & pTN_2 & pFP_3 & pTN_3 \\ \\ P(v_1) & & P(v_2) & & P(v_3) & \\ \\ \begin{array}{c} \sum_v P(v) pTP_v \\ \sum_v P(v) pFP_v \end{array} & & \begin{array}{c} \sum_v P(v) pFN_v \\ \sum_v P(v) pTN_v \end{array} & & & \\ \\ F = \frac{2pTP}{2pTP+pFN+pFP} \end{array}$$

In the ROC space

36 classifiers evaluated on the CDNET 2014 dataset in the ROC space:

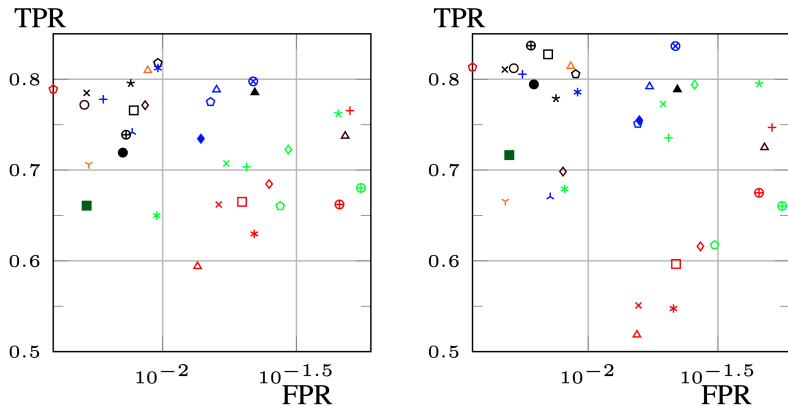
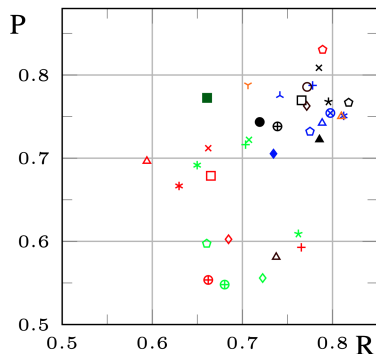
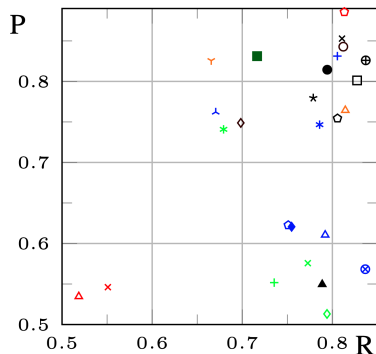


Figure: Summarized performances according to two different procedures in the cropped ROC space.

In the PR space



(c) CDNET procedure (PR space).



(d) Our procedure (PR space).

Figure: Summarized performances according to two different procedures in the cropped PR space.

Remember that **our summarization procedure preserves the bijection between the ROC and PR evaluation spaces!**

Ranking based on the F scores

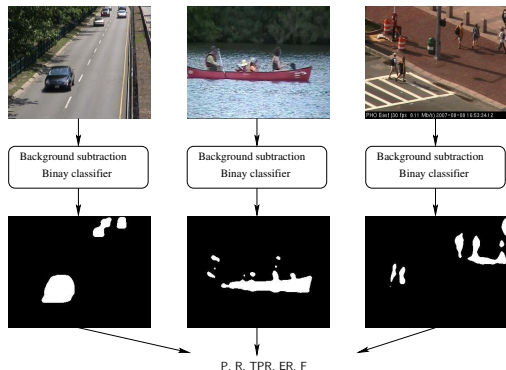
Algorithm	F of CDNET 2014	F [our summarization]
SemanticBGS	0.8098 (1)	0.8479 (1)
IUTIS-5	0.7821 (2)	0.8312 (3)
IUTIS-3	0.7694 (3)	0.8182 (5)
WisenetMD	0.7559 (4)	0.7791 (10)
SharedModel	0.7569 (5)	0.7885 (8)
WeSamBE	0.7491 (6)	0.7792 (9)
SuBSENSE	0.7453 (7)	0.7657 (12)
PAWCS	0.7478 (8)	0.8272 (4)

Table: Extract of F scores (and ranks) obtained with two procedures on CDNET.

Outline

- 1 Performance indicators for one video
- 2 Summarizing the performance for several videos
- 3 Summarizing applied on CDNET 2014
- 4 Conclusion

Take-home messages



- 1 It is unsound to average performance indicators, such as P, TPR, with the arithmetic mean because
 - 1 it breaks the consistency between indicators
 - 2 it makes the interpretation less reliable
- 2 **Prefer the summarization formulas**

More on summarization:

<http://www.telecom.ulg.ac.be/summarization>