



# Facing Imbalanced Data

## Recommendations for the Use of Performance Metrics



Laszlo A. Jeni <sup>1</sup>

Jeffrey F. Cohn <sup>1,2</sup>

Fernando De La Torre <sup>1</sup>

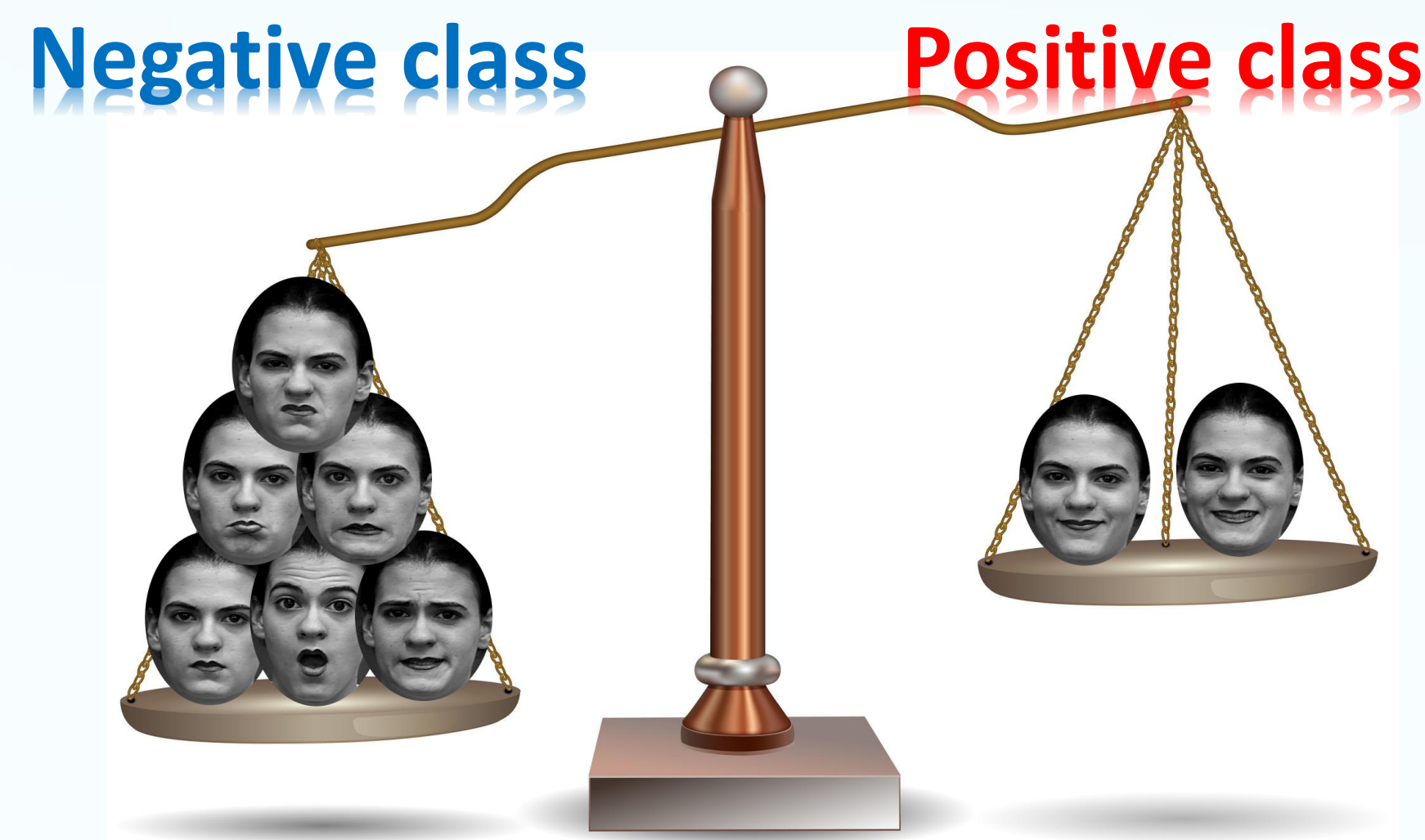
<sup>1</sup> Carnegie Mellon University, Pittsburgh, PA, [laszlo.jeni@ieee.org](mailto:laszlo.jeni@ieee.org), [ftorre@cs.cmu.edu](mailto:ftorre@cs.cmu.edu)

<sup>2</sup> University of Pittsburgh, Pittsburgh, PA, [jeffcohn@cs.cmu.edu](mailto:jeffcohn@cs.cmu.edu)

### 1. Introduction

Previous work on facial action unit (AU) recognition has emphasized face tracking and registration and the choice of features classifiers. Relatively little attention has been paid to **how imbalanced data may spoil performance metrics**.

**Facial expression data is typically highly skewed.** Imbalance in the test data distribution might produce misleading conclusions with certain metrics.



**Question: Is  $F_1 = 0.3$  good or bad performance?**

**Answer: It depends on the skew in the TEST set!**

### 3. Performance Metrics

In a binary classification problem the labels are either positive or negative. The decision made by the classifier can be represented as a 2x2 confusion matrix.

Predicted Class	Actual Class	
	1	-1
1	True Positives (TP)	False Positives (FP)
-1	False Negatives (FN)	True Negatives (TN)

#### Threshold Metrics:

**Accuracy** is the percentage of correctly classified positive and negative examples.

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

**Precision** is the fraction of recognized instances that are relevant.

$$\text{Prec} = \frac{TP}{TP+FP}$$

**Recall** is the fraction of relevant instances that are retrieved.

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

**$F_1$  score** is the harmonic mean of precision and recall.

$$F_1 = 2 \cdot \frac{\text{Prec} \cdot \text{Rec}}{\text{Prec} + \text{Rec}}$$

**Cohen's kappa** is observed agreement normalized to expected agreement.

$$K = \frac{P_{\text{Obs}} - P_{\text{Chance}}}{1 - P_{\text{Chance}}}$$

**Krippendorff's  $\alpha$**  is observed disagreement normalized to expected disagreement

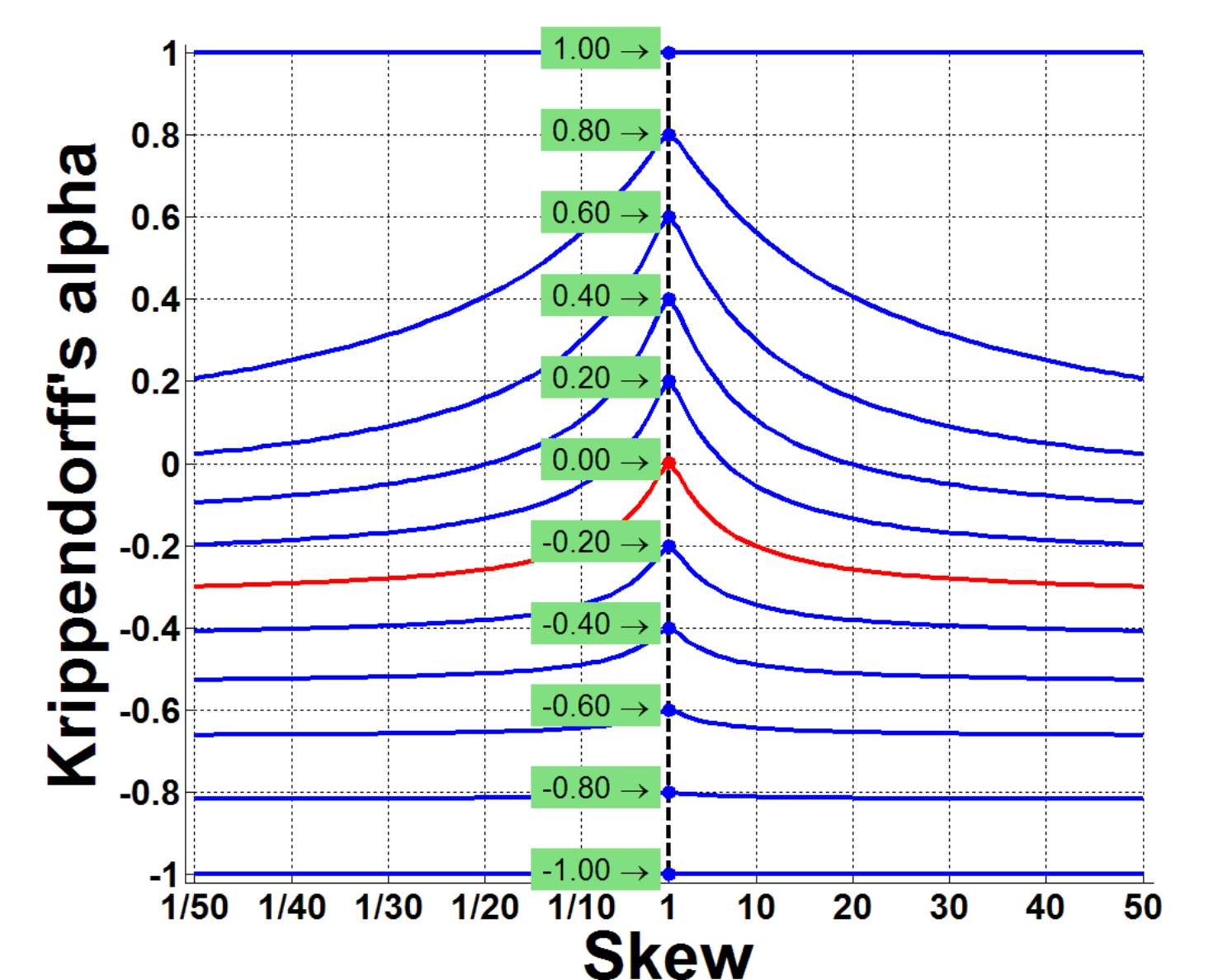
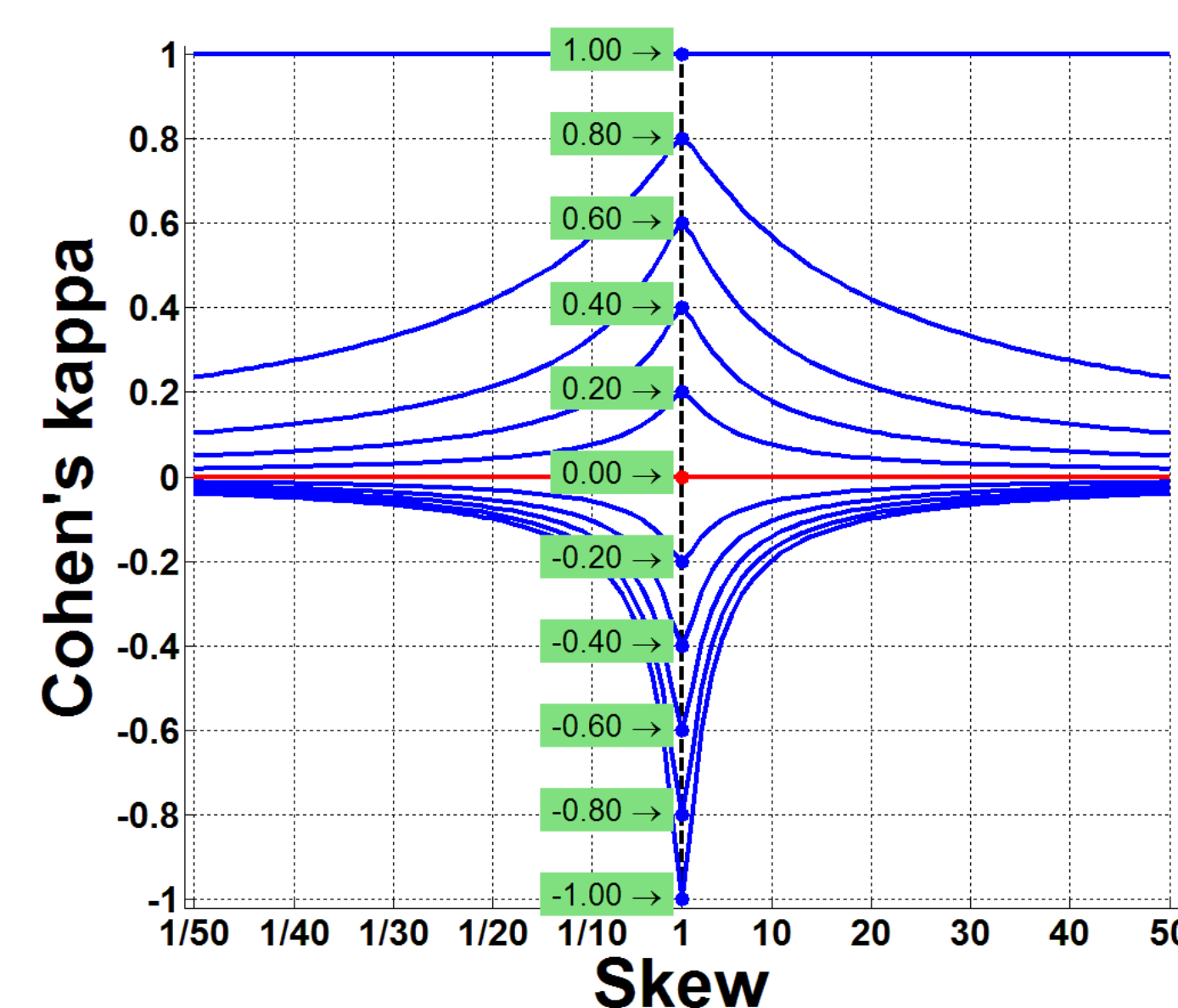
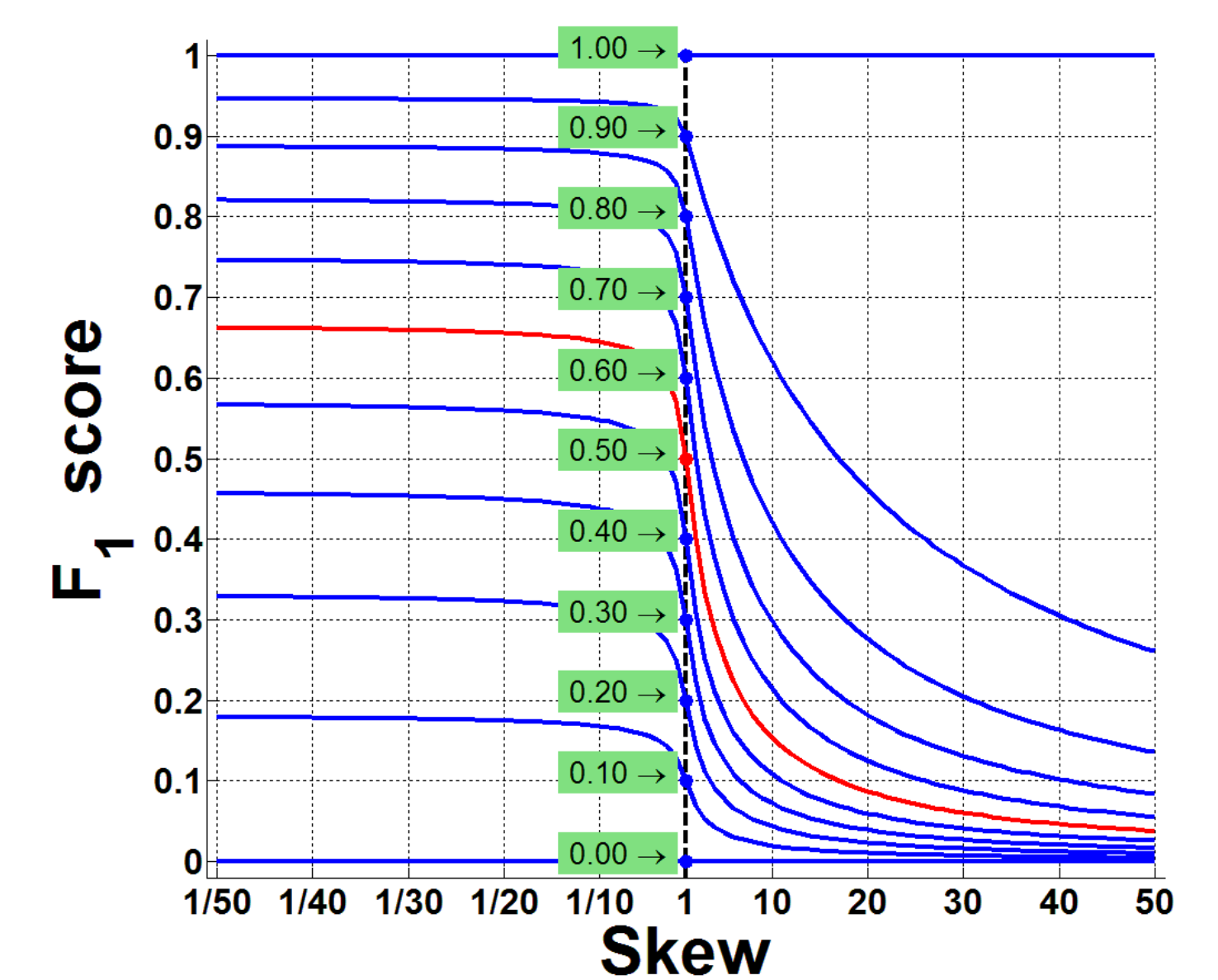
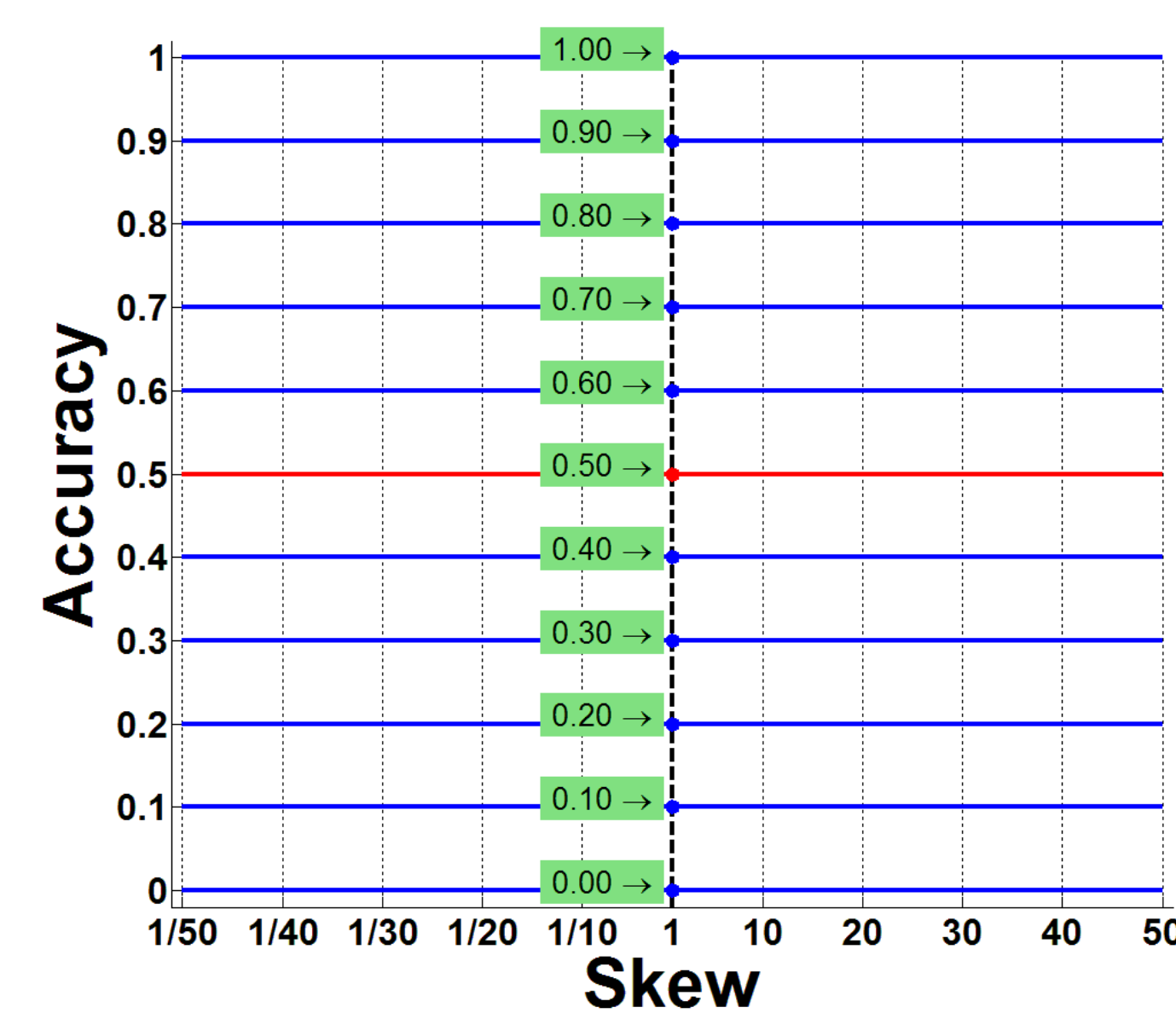
$$\alpha = 1 - \frac{D_{\text{Obs}}}{D_{\text{Chance}}}$$

#### Rank Metrics:

**ROC curve** shows the true positive rate as a function of false positive rate.

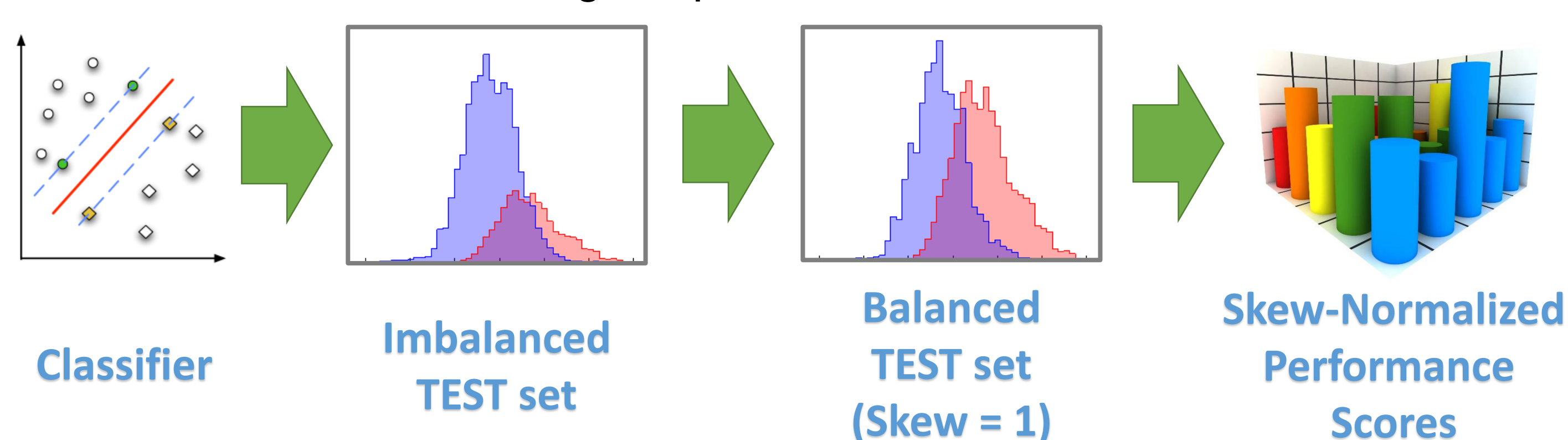
**Precision-Recall curve** shows the precision as the function of recall.

### 4. Behavior of Different Metrics



### 5. Skew Normalization

Different forms of re-sampling such as random over- and under-sampling can be used to balance the skewed distribution of the **TEST partitions** of the dataset before calculating the performance metrics.



Reporting both obtained performance metrics and skew-normalized scores, classifiers can be compared across databases free of confounds introduced by skew.

Code to compute skew-normalized scores for all of the metrics considered above and visualizations are available from:

<http://www.pitt.edu/~jeffcohn/skew/>

### 6. Results on Real Data

	AU	F <sub>1</sub>		Kappa		Alpha		AUC-ROC	
		Orig.	Norm.	Orig.	Norm.	Orig.	Norm.	Orig.	Norm.
CK+	4	0.73	0.83	0.68	0.69	0.68	0.68	0.90	0.90
	6	0.74	0.83	0.72	0.82	0.76	0.82	0.93	0.94
	9	0.92	0.97	0.92	0.96	0.92	0.96	1.00	1.00
	12	0.88	0.94	0.87	0.88	0.87	0.88	0.98	0.98
Pain Archive	4	0.06	0.67	0.11	0.38	0.11	0.38	0.74	0.75
	6	0.41	0.70	0.35	0.41	0.35	0.41	0.77	0.77
	9	0.20	0.69	0.20	0.48	0.20	0.47	0.75	0.75
	12	0.32	0.66	0.23	0.34	0.23	0.33	0.72	0.73
RU-FACS	4	0.00	0.52	0.01	0.10	0.00	0.06	0.55	0.53
	6	0.49	0.85	0.48	0.72	0.48	0.72	0.90	0.90
	9	0.00	0.68	0.00	0.47	0.00	0.47	0.71	0.68
	12	0.68	0.84	0.64	0.70	0.64	0.70	0.91	0.91

PERFORMANCE SCORES FOR THE ORIGINAL AND THE SKEW = 1 NORMALIZED VERSION OF THE THREE DATASETS (FOR MORE AUs, SEE TEXT).