#### Behavior of Distance-Based Methods in a Context of Class-Imbalance or High-Dimensionality

#### Rémi Emonet

Université Jean-Monnet, Laboratoire Hubert Curien, Saint-Étienne

Talk at Tahiti (Bréhat), 2019-06-27









# \$whoami

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### The Curse of Dimensionality

High-dimensionality is<sup>can be</sup> a mess.

#### What is this Curse Anyway?

• Some definition:

Various phenomena that arise when analyzing and organizing data in high-dimensional spaces.

- Term coined by Richard E. Bellman
  - 1920 1984
  - dynamic programming
  - differential equations
  - shortest path
- What is (not) the cause?
  - not an intrinsic property of the data
  - depends on the representation
  - depends on how data is analyzed

#### **Combinatorial Explosion**

- Suppose
  - $\circ$  you have d entities
  - each can be in 2 states
- Then
  - $\circ$  there are  $2^d$  combinations to consider/test/evaluate
- Happens when considering
  - $\circ$  all possible subsets of a set  $(2^d)$
  - $\circ$  all permutations of a list (d!)
  - $\circ$  all affectations of entities to labels ( $k^d$ , with k labels)

{ <b>a</b> }	{a,b}	{a,b,c}	{a,b,c,d}
{        } {a}	{	<pre>{      } {      c} {      b     } {      b,c} {a      } {a,      c} {a,b     } {a,b,c}</pre>	<pre>{</pre>
			ja, uj

## **Regular Space Coverage**

- Analogous to combinatorial explosion, in continuous spaces
- Happens when considering
  - histograms

Ο

- density estimation
- anomaly detection



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 8 / 56 (11/12)

## **In Modeling and Learning**

- The world is complicated
  - state with a huge number of variables (dimensions)
  - possibly noisy observations
  - e.g. a 1M-pixel image has 3 million dimensions
- Learning would need observations for each state
  - it would require too many examples
  - need for an "interpolation" procedure, to avoid overfitting
- Hughes phenomenon, 1968 paper (which is wrong, it seems)

given a (small) number of training samples,

additional feature measurements

may reduce the performance of a statistical classifier



#### A Focus on Distances/Volumes

- Considering a *d* dimensional space
- About volumes
  - $\circ$  volume of the cube:  $C_d(r) = (2r)^d$
  - $\circ$  volume of a sphere with radius r:  $S_d(r) = rac{\pi^{d/2}}{\Gamma(rac{d}{r}+1)}r^d$

 $(\Gamma ext{ is the continuous generalization of the factorial}) \circ ext{ ratio: } rac{S_d(r)}{C_d(r)} o 0 ext{ (linked to space coverage)}$ 



#### A Focus on Distances/Volumes (cont'd)





- average (euclidean) distance between two random points?
- everything becomes almost **as** "far"
- Happens when considering
  - radial distributions (multivariate normal, etc)
  - k-nearest neighbors (hubness problem)
  - other distance-based algorithms



#### The Curse of Dimensionality

Many things get degenerated with high dimensions Problem of: approach + data representation We have to hope that there is no curse

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### Ockham's Razor

#### Shave unnecessary assumptions.

dech 1 se .

ANA DE MARRIED MARRIED MAR

#### **Ockham's Razor**

- Term from 1852, in reference to Ockham (XIV<sup>th</sup>)
- *lex parsimoniae*, law of parsimony
- Prefer the simplest hypothesis that fits the data.
- Formulations by Ockham, but also earlier and later
- More a concept than a rule
  - simplicity
  - parsimony
  - elegance
  - shortness of explanation
  - shortness of program (Kolmogorov complexity)
  - falsifiability (sciencific method)
- According to Jürgen Schmidhuber, the appropriate mathematical theory of Occam's razor already exists, namely, Solomonoff's theory of optimal inductive inference.

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### Simplicity of Data: subspaces

- Data might be high-dimensional, but we have hope
  - that there is a organization or regularity in the highdimensionality
  - $\circ~$  that we can guess it
  - $\circ~$  or, that we can learn/find it
- Approaches: dimensionality reduction, manifold learning
   PCA, kPCA, \*PCA, SOM, Isomap, GPLVM, LLE, NMF, ...

#### Simplicity of Data: compressibility



- Idea
  - data can be high dimensional but compressible
  - i.e., there exist a compact representation
- Program that generates the data (Kolmogorov complexity)
- Sparse representations
  - wavelets (jpeg), fourier transform
  - sparse coding, representation learning
- Minimum description length
  - size of the "code" + size of the encoded data



#### **Simplicity of Models: information criteria**

- Used to select a model
- Penalizes by the number *k* of *free parameters* 
  - AIC (Aikake Information Criterion)
    - penalizes the Negative-Log-Likelihood by k
  - BIC (Bayesian IC)
    - penalizes the NLL by  $k \log(n)$  (for *n* observations)
  - BPIC (Bayesian Predictive IC)
  - DIC (Deviance IC)
  - FIC (Focused IC)
  - Hannan-Quinn IC
  - TIC (Takeuchi IC)
- Sparsity of the parameter vector (*l*0 norm)
  - penalizes the number of non-zero parameters

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- <u>High-dimensionality and Neighborhood</u>
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### A Focus on Distances/Volumes (cont'd)





- average (euclidean) distance between two random points?
- everything becomes almost **as** "far"
- Happens when considering
  - radial distributions (multivariate normal, etc)
  - k-nearest neighbors (hubness problem)
  - other distance-based algorithms



#### **Distance Contraction**

- Experiment
  - $\circ~$  sampling uniformly random points in the unit cube
  - looking at the distribution of inter-point distances
  - variance decreases with dimensionality



• Question: is it a problem? maybe not if the ranking is right

#### **Hubness Problem**

- Experiment
  - sampling uniformly random points in the unit cube
  - computing how often each point is in the nearest neighbor of another point  $\mathbb{E}[(N-\mu_N)^3]$
  - Hubness as skewness: *hubness* =



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 23 / 56

#### Hubness: why are some points so special?



hubness-uniform-sqeuclidean-False

- Where are these points?
- The border theory...
- ... so it is distribution-dependant

#### Hubness: testing the border theory

#### Wrapping the points (hyper-torus) lacksquare









200 400 600 800 1000

Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 25 / 56

#### Hubness: what is a border?

Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 26 / 56

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### **Imbalanced Problems: Examples**

- Anomaly detection
  - unsafe situations in videos
  - defect detection in images
  - abnormal heart beat detection in ECG
- Fraud detection
  - fraudulent checks
  - credit card fraud (physical, online)
  - financial fraud (French DGFIP)

#### **Imbalanced Classification Problems**

- Binary classification
  - + positive class: minority class, anomaly, rare event, ...
  - negative class: majority class, normality, typical event, ...
- Confusion matrix (of a model vs a ground truth)
  - TP: true positive
  - FP: false positive
  - TN: true negative
  - FN: false negative
- Some measures
  - Precision:  $prec = \frac{TP}{TP + FP}$
  - Recall:  $rec = \frac{TP}{P} = \frac{TP}{TP + FN}$
  - $\circ \ \ F_{eta}$ -measure:  $F_{eta} = (1+eta^2) rac{prec \cdot rec}{eta^2 \cdot prec + rec}$

\*(higher is better)

Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 29 / 56 (2/3)



#### F-measure vs Accuracy?

$$F_eta = (1+eta^2)rac{prec \cdot rec}{eta^2 \cdot prec + rec} = rac{(1+eta^2) \cdot (P-FN)}{1+eta^2 P - FN + FP}$$

$$accuracy = \frac{TP + TN}{P + N} = 1 - \frac{FN + FP}{P + N}$$

- Accuracy inadequacy (e.g. N = 10000, P = 100)
  - $\circ$  lazy "all-" classifier (TP = 0, TN = N, FP = 0, FN = P)

- $F_{\beta}$ -measure challenges
  - discrete (like the accuracy)
  - non-convex (even with continuous surrogates)
  - **non-separable**, i.e.  $F_{\beta} \neq \sum_{(x_i,y_i)\in S} ...$

#### Ok, but I'm doing gradient descent, so ...



- Gradient:  $0.2 \Rightarrow -7.21$ ,  $0.5 \Rightarrow -2.89$ ,  $0.8 \Rightarrow -1.80$ ,  $1 \Rightarrow -1.44$
- Example, gradient intensity is the same for:
  - $\circ~10+$  wrongly classified with an output proba. of 0.2
  - $\circ$  40 correctly classified with an output proba 0.8
  - $\circ$  i.e., lazily predicting systematically 0.2 (for +) yields a "stable" solution with 10+ vs 40-

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - <u>Reweight, resampling, etc</u>
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### **Counteracting Imbalance**

- Undersampling the majority class –
- Oversampling class +
- Generating fake +
- Using a weighted-classifiers learner

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

A Corrected Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data

 <u>Rémi Viola</u>, Rémi Emonet , Amaury Habrard, <u>Guillaume Metzler</u>, Sébastien Riou, Marc Sebban
 ???

#### k-NN: k Nearest Neighbor Classification

- k-NN
  - to classify a new point
  - find the closest k points (in the training section)
  - $\circ~$  use a voting scheme to affect a class
  - efficient algorithms (K-D Tree, Ball Tree)
- Does k-NN still matter?
  - non-linear by design (with similarity to RBF-kernel SVM)
  - no learning, easy to patch a model (add/remove points)
  - Limits of k-NN for imbalanced data?



#### Limits of k-NN for imbalanced data?

- 1. k-NN behavior in uncertain areas
  - $\circ~$  i.e., for some feature vector, the class can be + or -
  - $\circ$  i.e., the Bayes Risk is non zero
  - ✓ not so bad (respects imbalance)
- 2. k-NN behavior around boundaries
  - i.e., what happens if classes are separate but imbalanced
  - **X** sampling effects cause problems

#### k-NN at a boundary (1000 +)



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 38 / 56

## k-NN at a boundary (100 +)



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 39 / 56

## k-NN at a boundary (10 +)



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 40 / 56

#### k-NN: increasing k?



#### A Corrected Nearest Neighbor Algorithm Maximizing the F-Measure from Imbalanced Data

 Rémi Viola, Rémi Emonet, Amaury Habrard, Guillaume Metzler, Sébastien Riou, Marc Sebban
 ???

#### $\gamma$ -NN Idea: push the decision boundary



- Goal: correct for problems due to sampling with imbalance
- Genesis: GAN to generate "+" around existing ones
   ⇒ unstable, failing, complex
- Approach
  - artificially make + closer to new points
  - $\circ~$  how? by using a different distance for + and -
  - the base distance to + gets multiplied by a parameter  $\gamma$ (intuitively  $\gamma \leq 1$  if + is rare)

$$d_\gamma(x,x_i) = egin{cases} d(x,x_i) & ext{if } x_i \in S_-, \ \gamma \cdot d(x,x_i) & ext{if } x_i \in S_+. \end{cases}$$

Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 43 / 56 (3/4)

#### $\gamma$ -NN: varying $\gamma$ with two points



Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 44 / 56 (5/6)

#### $\gamma$ -NN: varying $\gamma$ with a few +



•  $\gamma$ -NN can control

how close to the minuses it pushes the boundary

#### $\gamma$ -NN: Algorithm

**Algorithm 1:** Classification of a new example with  $\gamma k$ -NN

**Input** : a query **x** to be classified, a set of labeled samples  $S = S_+ \cup S_-$ , a number of neighbors k, a positive real value  $\gamma$ , a distance function d**Output:** the predicted label of **x** 

 $\begin{array}{l} \mathcal{NN}^{-}, \mathcal{D}^{-} \leftarrow nn(k, \mathbf{x}, S_{-}) & // \text{ nearest negative neighbors with their distances} \\ \mathcal{NN}^{+}, \mathcal{D}^{+} \leftarrow nn(k, \mathbf{x}, S_{+}) & // \text{ nearest positive neighbors with their distances} \\ \mathcal{D}^{+} \leftarrow \gamma \cdot \mathcal{D}^{+} \\ \mathcal{NN}_{\gamma} \leftarrow firstK \left( k, sortedMerge((\mathcal{NN}^{-}, \mathcal{D}^{-}), (\mathcal{NN}^{+}, \mathcal{D}^{+})) \right) \\ y \leftarrow + \text{ if } \left| \mathcal{NN}_{\gamma} \cap \mathcal{NN}^{+} \right| \geq \frac{k}{2} \text{ else } - // \text{ majority vote based on } \mathcal{NN}_{\gamma} \\ \textbf{return } y \end{array}$ 

- Trivial to implement
- Same complexity as k-NN (at most twice)
- Training
  - ∘ none, as k-NN
  - γ is selected by cross-validation (on the measure of interest)

#### $\gamma\text{-NN}:$ a way to reweight distributions

- In uncertain regions
- At the boundaries

## **Results on public datasets (F-measure)**

DATASETS	3-NN	DUPk-NN	wk-NN	CWk-NN	LMNN	$\gamma k$ -NN
BALANCE	0.954(0.017)	0.954(0.017)	$0.957 \scriptscriptstyle (0.017)$	$0.961 \scriptscriptstyle (0.010)$	$0.963_{(0.012)}$	$0.954 \scriptscriptstyle (0.029)$
AUTOMPG	$0.808_{(0.077)}$	$0.826 \scriptscriptstyle (0.033)$	0.810(0.076)	$0.815 \scriptscriptstyle (0.053)$	$0.827 \scriptscriptstyle (0.054)$	$0.831_{(0.025)}$
IONO	$0.752 \scriptscriptstyle (0.053)$	$0.859 \scriptscriptstyle (0.021)$	$0.756 \scriptscriptstyle (0.060)$	$0.799 \scriptscriptstyle (0.036)$	$0.890 \scriptscriptstyle (0.039)$	$0.925_{(0.017)}$
PIMA	0.500(0.056)	$0.539 \scriptscriptstyle (0.033)$	$0.479 \scriptscriptstyle (0.044)$	$0.515 \scriptscriptstyle (0.037)$	$0.499 \scriptscriptstyle (0.070)$	0.560(0.024)
WINE	$0.881_{(0.072)}$	0.852(0.057)	0.881(0.072)	$0.876 \scriptscriptstyle (0.080)$	0.950(0.036)	$0.856 \scriptscriptstyle (0.086)$
GLASS	$0.727 \scriptscriptstyle (0.049)$	$0.733 \scriptscriptstyle (0.061)$	$0.736 \scriptscriptstyle (0.052)$	$0.717 \scriptscriptstyle (0.055)$	$0.725 \scriptscriptstyle (0.048)$	$0.746_{(0.046)}$
GERMAN	$0.330 \scriptscriptstyle (0.030)$	0.449(0.037)	$0.326 \scriptscriptstyle (0.030)$	$0.344 \scriptscriptstyle (0.029)$	$0.323 \scriptscriptstyle (0.054)$	0.464 (0.029)
VEHICLE	$0.891_{(0.044)}$	0.867(0.027)	0.891(0.044)	$0.881 \scriptscriptstyle (0.021)$	$0.958_{(0.020)}$	0.880(0.049)
HAYES	$0.036 \scriptscriptstyle (0.081)$	0.183(0.130)	0.050(0.112)	$0.221 \scriptscriptstyle (0.133)$	$0.036 \scriptscriptstyle (0.081)$	0.593 (0.072)
SEGMENTATION	$0.859 \scriptscriptstyle (0.028)$	0.862(0.018)	0.877 (0.028)	$0.851 \scriptscriptstyle (0.022)$	$0.885_{(0.034)}$	$0.848 \scriptscriptstyle (0.025)$
ABALONE8	$0.243 \scriptscriptstyle (0.037)$	$0.318 \scriptscriptstyle (0.013)$	$0.241 \scriptscriptstyle (0.034)$	$0.330 \scriptscriptstyle (0.015)$	$0.246 \scriptscriptstyle (0.065)$	$0.349_{(0.018)}$
yeast3	$0.634 \scriptscriptstyle (0.066)$	$0.670 \scriptscriptstyle (0.034)$	$0.634 \scriptscriptstyle (0.066)$	$0.699_{(0.015)}$	$0.667 \scriptscriptstyle (0.055)$	0.687 (0.033)
PAGEBLOCKS	$0.842 \scriptscriptstyle (0.020)$	$0.850 \scriptscriptstyle (0.024)$	$0.849 \scriptscriptstyle (0.019)$	$0.847 \scriptscriptstyle (0.029)$	$0.856_{(0.032)}$	$0.844 \scriptscriptstyle (0.023)$
SATIMAGE	$0.454 \scriptscriptstyle (0.039)$	0.457(0.027)	0.454(0.039)	$0.457 \scriptscriptstyle (0.023)$	$0.487_{(0.026)}$	0.430(0.008)
LIBRAS	$0.806_{(0.076)}$	$0.788_{(0.187)}$	$0.806_{(0.076)}$	$0.789 \scriptscriptstyle (0.097)$	0.770(0.027)	$0.768 \scriptscriptstyle (0.106)$
WINE4	$0.031_{(0.069)}$	0.090(0.086)	0.031 (0.069)	$0.019 \scriptscriptstyle (0.042)$	0.000(0.000)	0.090(0.036)
yeast6	$0.503 \scriptscriptstyle (0.302)$	$0.449_{(0.112)}$	$0.502 \scriptscriptstyle (0.297)$	$0.338 \scriptscriptstyle (0.071)$	$0.505 \scriptscriptstyle (0.231)$	$0.553_{(0.215)}$
ABALONE17	$0.057_{(0.078)}$	$0.172_{(0.086)}$	0.057 (0.078)	$0.096 \scriptscriptstyle (0.059)$	0.000(0.000)	0.100(0.038)
ABALONE20	$0.000_{(0.000)}$	0.000(0.000)	0.000(0.000)	$0.067_{(0.038)}$	$0.057 \scriptscriptstyle (0.128)$	$0.052 \scriptscriptstyle (0.047)$
MEAN	$0.543 \scriptscriptstyle (0.063)$	$0.575 \scriptscriptstyle (0.053)$	$0.544 \scriptscriptstyle (0.064)$	$0.559 \scriptscriptstyle (0.046)$	$0.560 \scriptscriptstyle (0.053)$	0.607(0.049)

#### **Results on DGFiP datasets (F-measure)**

DATASETS	3-NN	$\gamma k - NN$	SMOTE	$ SMOTE + \gamma k - NN $
Dgfip19 2	$0,\!454 \scriptscriptstyle (0,007)$	$0,528 \scriptscriptstyle (0,005)$	$0,505_{(0,010)}$	$0,529_{(0,003)}$
Dgfip9 2	$0,\!173 \scriptscriptstyle (0,074)$	$\overline{0,\!396}_{(0,018)}$	$0,340_{(0,033)}$	$0,419_{(0,029)}$
DGFIP $4\ 2$	$0,\!164 \scriptscriptstyle (0,155)$	$\overline{0,\!373}_{(0,018)}$	$0,368_{(0,057)}$	$0,377_{(0,018)}$
Dgfip8 1	$0,100_{(0,045)}$	$\overline{0,299}_{(0,010)}$	$0,278_{(0,043)}$	$0,299_{(0,011)}$
Dgfip8 2	$0,140_{(0,078)}$	$0,\!292_{(0,028)}$	0,313(0,048)	$0,312 \scriptscriptstyle (0,021)$
Dgfip9 1	$0,088 \scriptscriptstyle (0,090)$	$0,\!258_{(0,036)}$	$0,270_{(0,079)}$	$\overline{0,288}_{(0,026)}$
Dgfip4 1	$0,\!073_{(0,101)}$	$0,\!231_{(0,139)}$	$\overline{0,199}_{(0,129)}$	$0,278_{(0,067)}$
Dgfip16 1	$0,049_{(0,074)}$	$\overline{0,166}_{(0,065)}$	$0,180_{(0,061)}$	$0,191_{(0,081)}$
Dgfip16 2	$0,210_{(0,102)}$	$0,\!202_{(0,056)}$	$\overline{0,220}_{(0,043)}$	$0,229_{(0,026)}$
Dgfip20 3	$0,142 \scriptscriptstyle (0,015)$	$0,210_{(0,019)}$	$\overline{0,199}_{(0,015)}$	$0,212_{(0,019)}$
Dgfip5 3	$0,\!030 \scriptscriptstyle (0,012)$	$\overline{0,\!105}_{(0,008)}$	$0,110_{(0,109)}$	$\underline{0,107}_{(0,010)}$
MEAN	$0,148 \scriptscriptstyle (0,068)$	$\boxed{0,\!278}_{\scriptscriptstyle (0,037)}$	$0,\!271_{(0,057)}$	<b>0,295</b> (0,028)

# $\gamma$ -NN at a boundary (10 and 100 +)



#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- Discussion

#### From Cost-Sensitive Classification to Tight F-measure Bounds

- <u>Kevin Bascol</u>, Rémi Emonet, Elisa Fromont, Amaury Habrard, <u>Guillaume Metzler</u>, Marc Sebban
- AISTATS2019

#### **Optimizing the** $F_{\beta}$ **-measure?**

• Reminder

• Precision: 
$$prec = \frac{TP}{TP + FP}$$

• Recall: 
$$rec = \frac{TP}{P} = \frac{TP}{TP + FN}$$

$$\circ \ F_{eta}$$
-measure:  $F_{eta} = (1+eta^2) rac{prec \cdot rec}{eta^2 \cdot prec + rec}$ 

• Non-separability, i.e.  $F_{\beta} \neq \sum_{(x_i,y_i)\in S} ...$ 

NB: accuracy is separable,  $acc = \sum_{(x_i,y_i)\in S} rac{1}{m} \delta(y_i - \hat{y_i})$ 

⇒ The loss for one point depends on the others
 ⇒ Impossible to optimize directly
 ⇒ Impossible to optimize on a subset (minibatch)

#### Weighted classification for $F_{\beta}$

$$F_{\beta} = \frac{(1+\beta^2) \cdot (P-FN)}{1+\beta^2 P - FN + FP} = \frac{(1+\beta^2) \cdot (P-e_1)}{1+\beta^2 P - e_1 + e_2}$$

- The  $F_{\beta}$ -measure is linear fractional (in  $e = (e_1, e_2) = (FN, FP)$ ) i.e.  $F_{\beta} = \frac{\langle a', e \rangle + b}{\langle c, e \rangle + d} = \frac{A}{B}$
- Relation to weighted classification

 $F_{\beta} \geq t \quad (\text{we achieve a good, above } t, F_{\beta} \text{ value}) \\ \Leftrightarrow A \geq t \cdot B \\ \Leftrightarrow A - t \cdot B \geq 0 \\ \Leftrightarrow (1 + \beta^2) \cdot (P - e_1) - t(1 + \beta^2 P - e_1 + e_2) \geq 0 \\ \Leftrightarrow (-1 - \beta^2 + t)e_1 - te_2 \geq -P(1 + \beta^2) + t(1 + \beta^2 P) \\ \Leftrightarrow (1 + \beta^2 - t)e_1 + te_2 \leq -P(1 + \beta^2) + t(1 + \beta^2 P) \\ \Rightarrow \text{ so, we can minimize the weighted problem} \\ \text{ with class weights } a(t) = (1 + \beta^2 - t, t) \end{cases}$ 

#### **Overview**

- Introduction
- High-dimensional problems
  - The curse of dimensionality
  - Ockham's Razor
  - Notions of Simplicity
- High-dimensionality and Neighborhood
- Imbalanced classification problems
  - The Problem (and performance measures)
  - Reweight, resampling, etc
  - Correcting k-NN ( $\gamma$ -NN)
  - Focusing on the F-Measure optimization (Élisa)
- <u>Discussion</u>

#### Thank you! Questions?

and now for something completely different...

Tahiti (Bréhat) | Rémi Emonet | 2019-06-27 | 56 / 56