









### $\gamma$ -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**   $\mathcal{NN}^-, \mathcal{D}^- \leftarrow nn(k, \mathbf{x}, S_-)$  // nearest negative neighbors with their distances  $\mathcal{NN}^+, \mathcal{D}^+ \leftarrow nn(k, \mathbf{x}, S_+)$  // nearest positive neighbors with their distances  $\mathcal{D}^+ \leftarrow \gamma \cdot \mathcal{D}^+$ 

 $\mathcal{NN}_{\gamma} \leftarrow firstK\left(k, sortedMerge((\mathcal{NN}^{-}, \mathcal{D}^{-}), (\mathcal{NN}^{+}, \mathcal{D}^{+}))\right)$ 

 $\begin{array}{ll} y \leftarrow + \mbox{ if } \left| \mathcal{N} \mathcal{N}_{\gamma} \cap \mathcal{N} \mathcal{N}^+ \right| \geq \frac{k}{2} \mbox{ else } - & // \mbox{ majority vote based on } \mathcal{N} \mathcal{N}_{\gamma} \\ \mbox{ return } y \end{array}$ 

- Trivial to implement
- Same complexity as k-NN (at most twice)
- Training
  - $\circ~$  none, as k-NN  $\circ~\gamma$  is selected by (cross-)validation
    - (on the measure of interest)

SLEIGHT Science Event#6 | Rémi Emonet | 2021-07-06 | 31 / 92 (3/3

# Results on public datasets (F-measure) DATASETS 3-NN DUPk-NN wk-NN CWk-NN LMNN γk-NN BALANCE 0.954(0.017) 0.954(0.017) 0.957(0.017) 0.961(0.010) 0.963(0.012) 0.954(0.029)

AUTOMPG	0.808(0.077)	0.826(0.033)	0.810(0.076)	0.815(0.053)	0.827(0.054)	0.831(0.025)
IONO	0.752(0.653)	0.859(0.021)	0.756(0.060)	0.799(0.036)	0.890(0.639)	0.925(0.017)
PIMA	0.500(0.056)	0.539(0.033)	0.479(0.044)	0.515(0.037)	0.499(0.670)	0.560(0.024)
WINE	0.881(0.072)	0.852(0.057)	0.881(0.072)	0.876(0.050)	0.950(0.036)	0.856(0.056)
GLASS	0.727(0.049)	0.733(0.001)	0.736(0.052)	0.717(0.055)	0.725(0.048)	0.746(0.040)
GERMAN	0.330(0.030)	0.449(0.037)	0.326(0.030)	0.344(0.029)	0.323(0.054)	0.464(0.029)
VEHICLE	$0.891_{(0.044)}$	0.867(0.027)	0.891(0.944)	0.881(0.021)	0.958(0.020)	0.880(0.049)
HAYES	0.036(0.081)	$0.183_{(0.130)}$	0.050(0.112)	0.221(0.133)	0.036(0.081)	0.593(0.072)
SEGMENTATION	0.859(0.028)	0.862(0.018)	0.877(0.028)	0.851(0.022)	0.885(0.034)	0.848(0.025)
ABALONE8	0.243(0.037)	0.318(0.013)	0.241(0.034)	0.330(0.015)	0.246(0.065)	0.349(0.018)
YEAST3	0.634(0.665)	0.670(0.034)	0.634(0.966)	0.699(0.015)	0.667(0.055)	0.687(0.033)
PAGEBLOCKS	0.842 (0.020)	0.850(0.024)	0.849(0.019)	0.847(0.029)	0.856(0.032)	0.844(0.023)
SATIMAGE	$0.454_{(0.039)}$	$0.457_{(0.027)}$	0.454(0.039)	0.457(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(0.057)	0.770(0.027)	0.768(0.106)
WINEA	0.031(0.069)	0.090(0.086)	0.031(0.069)	0.019(0.042)	0.000(0.000)	0.090(0.036)
YEAST6	0.503(0.302)	0.449(0.112)	0.502(0.297)	0.338(0.071)	0.505(0.231)	0.553(0.215)
ABALONE17	$0.057_{(0.078)}$	0.172(0.086)	0.057(0.078)	0.096(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(0.128)	$0.052 \scriptscriptstyle (0.047)$
MEAN	0.543 (1 1 1 2 2	0.5750000	0.544	0.559.000	0.560	0 607 (0 040)

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

- < In uncertain regions (1-NN is already ok)
- ✓ At the **boundaries** (10 and 100 +)





Results on DGFiP datasets (F-measure)						
DATASETS	3–NN	$\gamma k - \mathrm{NN}$	SMOTE	$SMOTE + \gamma k - NN$		
DGFIP $19\ 2$	$0,\!454 \scriptscriptstyle (0,007)$	$0,\!528 \scriptscriptstyle (0,005)$	$0,505\scriptscriptstyle (0,010)$	0,529(0,003)		
Dgfip9 2	$0,\!173 \scriptscriptstyle (0,074)$	$\overline{0,\!396}_{(0,018)}$	$0,\!340 \scriptscriptstyle (0,033)$	$0,419_{(0,029)}$		
DGFIP $4\ 2$	$0,\!164\scriptscriptstyle (0,155)$	$\overline{0,\!373}_{(0,018)}$	$0,368 \scriptscriptstyle (0,057)$	$0,377_{(0,018)}$		
Dgfip8 1	$0,100 \scriptscriptstyle (0,045)$	$\overline{0,299}_{(0,010)}$	$0,\!278 \scriptscriptstyle (0,043)$	$0,299_{(0,011)}$		
Dgfip8 2	$0,140 \scriptscriptstyle (0,078)$	$0,\!292_{(0,028)}$	0,313 (0,048)	$0,312_{(0,021)}$		
Dgfip9 1	$0,\!088 \scriptscriptstyle (0,090)$	$0,\!258_{(0,036)}$	$0,270_{(0,079)}$	$\overline{0,288}_{(0,026)}$		
Dgfip4 1	$0,073 \scriptscriptstyle (0,101)$	$0,\!231_{(0,139)}$	$\overline{0,199}_{\scriptscriptstyle (0,129)}$	$0,278_{(0,067)}$		
Dgfip $16\ 1$	$0,049 \scriptscriptstyle (0,074)$	$\overline{0,166}_{(0,065)}$	$0,\!180 \scriptscriptstyle (0,061)$	$0,191_{(0,081)}$		
DGFIP $16\ 2$	$0,210\scriptscriptstyle (0,102)$	$0,\!202_{(0,056)}$	$\overline{0,220}{\scriptstyle (0,043)}$	$0,229_{(0,026)}$		
Dgfip $20$ 3	$0,\!142 \scriptscriptstyle (0,015)$	$0,\!210_{(0,019)}$	$\overline{0,199}_{\scriptscriptstyle (0,015)}$	$0,212_{(0,019)}$		
Dgfip5 3	$0,\!030 \scriptscriptstyle (0,012)$	$\overline{0,\!105}_{(0,008)}$	$0,110_{(0,109)}$	$\underline{0,107}_{(0,010)}$		
MEAN	$0,148_{(0,068)}$	$\underline{0,\!278}_{\scriptscriptstyle (0,037)}$	$0,\!271 \scriptscriptstyle (0,057)$	$0,295_{(0,028)}$		

## Extension With Metric Learning

- Note:
  - $\circ \gamma$ -NN learns a metric for comparing a query to a +
  - $\circ \, \gamma$ -NN kind of learn the size of a sphere around +
  - $\circ~$  this is a very simple "Metric Learning"
- Extension
  - learn a full metric (a matrix *M* and not only a scalar γ)
    derive a learning algorithm (not just using a validation set)

Learning from Few Positives: a Provably Accurate Metric Learning Algorithm to deal with Imbalanced Data

- <u>Rémi Viola</u>, Rémi Emonet, Amaury Habrard, <u>Guillaume Metzler</u>, Marc Sebban
- IJCAI 2020 (International Joint Conference on Artificial Intelligence)  $\circ$  learn a full metric (a matrix *M* and not only a scalar  $\gamma$ )
  - derive a learning algorithm (not just using a validation set)
  - derive a theoretical guarantees

SLEIGHT Science Event#6 | Rémi Emonet | 2021-07-06 | 35 / 92 (2/2)

SLEIGHT Science Event#6 | Rémi Emonet | 2021-07-06 | 36 / 92











Simple (but realistic) unsupervised	<b>Process behind the</b> $3\sigma$ rule
anomaly detection	
<ul> <li>Setup and approach: the three-sigma rule <ul> <li>we have a set of (unlabeled) points</li> <li>we consider one feature of interest</li> <li>we look at the standard deviation σ and mean of this feature</li> <li>anything beyond 3σ is an outlier/anomaly</li> </ul> </li> <li>Possible improvements <ul> <li>use robust statistics (percentiles, robust estimation)</li> <li>use several features</li> </ul> </li> </ul>	<ul> <li>Overall process</li> <li>we suppose a parametric model that explains the data</li> <li>i.e., we suppose the data is generated by this model</li> <li>here, the data comes from the a normal distribution (with two unknown parameters: a mean μ and a standard deviation σ)</li> <li>we estimate the parameters from the data</li> <li>here, using the empirical mean and stdev</li> <li>if the likelihood of a (new) point is low, it is an outlier</li> <li>here, the normal density is very low after 3σ</li> </ul>
SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   67 / 92 (272)	SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   68 / 92 (2/2)
Overview	A few methods (ex. from scikit-learn)
<ul> <li>Introduction</li> <li>Anomaly and (rare) event detection <ul> <li>Problem, notations and performance measures</li> </ul> </li> <li>Imbalanced classification problems <ul> <li>General approaches</li> <li>Correcting k-NN: γ-NN and MLFP</li> <li>Learning Maximum Excluding Ellipsoids</li> <li>Focusing on the F-Measure optimization</li> </ul> </li> <li>Unsupervised anomaly detection <ul> <li>Simple motivating approach</li> <li>A variety of ML methods</li> <li>Probabilistic models 101</li> <li>Case study: temporal motif mining</li> <li>More models (VAE, GAN, tensor networks,)</li> </ul> </li> <li>Closing remarks</li> </ul>	
SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   69 / 92	SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   70 / 92
Overview	Probabilistic Generative Models
<ul> <li>Introduction</li> <li>Anomaly and (rare) event detection <ul> <li>Problem, notations and performance measures</li> </ul> </li> <li>Imbalanced classification problems <ul> <li>General approaches</li> <li>Correcting k-NN: γ-NN and MLFP</li> <li>Learning Maximum Excluding Ellipsoids</li> <li>Focusing on the F-Measure optimization</li> </ul> </li> <li>Unsupervised anomaly detection <ul> <li>Simple motivating approach</li> <li>A variety of ML methods</li> <li><u>Probabilistic models 101</u></li> <li>Case study: temporal motif mining</li> <li>More models (VAE, GAN, tensor networks,)</li> </ul> </li> <li>Closing remarks</li> </ul>	<ul> <li>A generalized "three sigma" rule <ul> <li>more complex models</li> <li>more application than just anomaly detection</li> </ul> </li> <li>Modeling step <ul> <li>we define a (stochastic) generative story</li> <li>we define how we suppose data are generated</li> <li>we encode our knowledge/assumptions/constraints</li> <li>we define what is random and what is a parameter</li> </ul> </li> <li>Learning/fitting step <ul> <li>given the data, what do we know about parameters</li> <li>given the data, we find <i>the best</i> parameters</li> </ul> </li> <li>Some possible usage <ul> <li>given the learned parameters, what is an outlier?</li> <li>what do the parameters look like?</li> <li>what data could I generate from these parameters?</li> </ul> </li> </ul>
SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   71 / 92	SLEIGHT Science Event#6   Rémi Emonet   2021-07-06   72 / 92 (4/4)







