Extreme Multi-label Classification for Information Retrieval XMLC4IR Tutorial at ECIR 2018

Rohit Babbar¹ and Krzysztof Dembczyński²

¹Aalto University, Finland, ²Poznań University of Technology, Poland



European Conference on Information Retrieval Grenoble, France, March 26, 2018

• Rohit Babbar:

- Affiliation: University Aalto
- Previous affiliations: Max-Planck Institute for Intelligent Systems, Université Grenoble Alpes
- Main interests: extreme classification, multi-label classification, multi-class classification, text classification



- Krzysztof Dembczyński:
 - Affiliation: Poznan University of Technology
 - Previous affiliations: Marburg University
 - Main interests: extreme multi-machine learning, multi-label classification, label tree algorithms, learning theory



Agenda

- 1 Extreme classification: applications and challenges
- Algorithms
 - Label embeddings
 - Smart 1-vs-All approaches
 - Tree-based methods
 - Label filtering/maximum inner product search
- 3 Live demonstration

Webpage: http://www.cs.put.poznan.pl/ kdembczynski/xmlc-tutorial-ecir-2018/

Agenda

1 Extreme classification: applications and challenges

- Algorithms
 - Label embeddings
 - Smart 1-vs-All approaches
 - Tree-based methods
 - Label filtering/maximum inner product search
- 3 Live demonstration

Webpage: http://www.cs.put.poznan.pl/ kdembczynski/xmlc-tutorial-ecir-2018/ **Extreme multi-label classification** is a problem of **labeling** an item with a **small** set of tags out of an **extremely large** number of potential tags.











Alan Turing, 1912 births, 1954 deaths 20th-century mathematicians, 20th-century philosophers Academics of the University of Manchester Institute of Science and Technology Alumni of King's College, Cambridge Artificial intelligence researchers Atheist philosophers, Bayesian statisticians, British cryptographers, British logicians British long-distance runners, British male athletes, British people of World War II Computability theorists, Computer designers, English atheists English computer scientists. English inventors. English logicians English long-distance runners, English mathematicians English people of Scottish descent, English philosophers, Former Protestants Fellows of the Royal Society. Gav men Government Communications Headquarters people, History of artificial intelligence Inventors who committed suicide, LGBT scientists LGBT scientists from the United Kingdom, Male long-distance runners Mathematicians who committed suicide. Officers of the Order of the British Empire People associated with Bletchley Park, People educated at Sherborne School People from Maida Vale, People from Wilmslow People prosecuted under anti-homosexuality laws. Philosophers of mind Philosophers who committed suicide. Princeton University alumni, 1930-39 Programmers who committed suicide, People who have received posthumous pardons Recipients of British royal pardons. Academics of the University of Manchester Suicides by cyanide poisoning. Suicides in England, Theoretical computer scientists



New question \Rightarrow Assignment/recommendation of users



Customers who viewed Convex Optimization also viewed: Introduction to Linear Optimization (Athena Scientific Series in Numerical Optimization (Springer Series in Operations Research Combinatorial Optimization: Algorithms and Complexity (Dover and FL. Books on... Opt... Buy new: \$69.18 Buy new: \$13.94 Buy new: \$79.43 88 Used & new from \$42.86 83 Used & new from \$7.51 42 Used & new from \$58.00 ★★★★★ (23) √Prime ★★★★★ (25) √Prime ***** (34) -Prime Books > Science & Math > Mathematics Convex Optimization 1st Edition by Stephen Boyd (Author), Lieven Vandenberghe (Author) ***** 33 customer reviews Look inside J eTextbook 🖵 🗆 🛛 Hardcover Other Sellers \$63.12 \$33.08 - \$82.85 from \$66.99 O Rent \$33.08 - \$33.09 Buy new \$82.86 Convex In Stock. List Price: \$99-09 Save: \$17.13 (17%) Optimization Ships from and sold by Amazon.com. Gift-wrap available. 40 New from \$67.58 FREE Shipping. Want it tomorrow, July 12? Order within 7 hrs 43 mins and choose One-Day Shipping at checkout. Details Qty: 1.0 ų Add to Cart Turn on 1-Click ordering Ship to: Select a shirping address: + ISBN-13: 978-0521833783 ISBN-10: 0521833787 More Buying Choices 70 used & new from \$66.99 Why is ISBN important? . 40 New from \$67.58 30 Used from \$66.99 See All Busing Options Trade in your item Trada in

Selected item \Rightarrow Recommendation of top 3 items



Sequence of words \Rightarrow Recommendation of the next word

ZURICH[®]



· Specialist 24 hour roadside assistance (available in the Republic of Ireland and Northern

Our electric car insurance cover includes

· 20% discount off your electric car insurance premium

Possible bid phrases:

- Zurich car insurance
- Car insurance
- Auto insurance
- Vehicle insurance
- Electric car insurance

On-line ad \Rightarrow Recommendation of queries to an advertiser

053 915 7775

1890 400 300



Suggestion of top Twitter Trends

Setting

• Multi-class classification:

$$\boldsymbol{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d \xrightarrow{h(\boldsymbol{x})} y \in \{1, \dots, m\}$$

| | x_1 | x_2 | x_d | y |
|------------------|-------|-------|-----------|---|
| \boldsymbol{x} | 4.0 | 2.5 | -1.5 | 5 |

Setting

| $\boldsymbol{x} = (x_1, x_2)$ | $x_2,$ | \ldots, x_d | $) \in \mathbb{R}^{d}$ | $h(\boldsymbol{x})$ | $\rightarrow y \in \mathcal{Y}$ | $\{1, \ldots,$ | , m |
|-------------------------------|----------------|---------------|------------------------|---------------------|---------------------------------|----------------|-----|
| | | x_1 | x_2 | | x_d | y | |
| | \overline{x} | 4.0 | 2.5 | | -1.5 | 5 | |

• Multi-class classification:

• Multi-label classification:

$$\boldsymbol{x} = (x_1, x_2, \dots, x_d) \in \mathbb{R}^d \xrightarrow{\boldsymbol{h}(\boldsymbol{x})} \boldsymbol{y} = (y_1, y_2, \dots, y_m) \in \{0, 1\}^m$$

| | x_1 | x_2 | x_d | y_1 | y_2 | y_m |
|------------------|-------|-------|-----------|-------|-------|-----------|
| \boldsymbol{x} | 4.0 | 2.5 | -1.5 | 1 | 1 | 0 |

}

Extreme classification \Rightarrow a large number of labels $m (\geq 10^5)$

• Predictive performance:

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - Learning theory for large m

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - Learning theory for large m
 - ► Training and prediction under limited time and space budged

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - \blacktriangleright Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - \blacktriangleright Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning
 - Long-tail label distributions and zero-shot learning

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - \blacktriangleright Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning
 - Long-tail label distributions and zero-shot learning
- Computational complexity:

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - \blacktriangleright Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning
 - Long-tail label distributions and zero-shot learning
- Computational complexity:
 - ► time vs. space

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - \blacktriangleright Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning
 - Long-tail label distributions and zero-shot learning
- Computational complexity:
 - ► time vs. space
 - ► #examples vs. #features vs. #labels

Extreme classification \Rightarrow a large number of labels $m \ (\geq 10^5)$

- Predictive performance:
 - ▶ Performance measures: Hamming loss, prec@k, NDCG@k, Macro F
 - Learning theory for large m
 - Training and prediction under limited time and space budged
 - ► Learning with missing labels and positive-unlabeled learning
 - Long-tail label distributions and zero-shot learning

• Computational complexity:

- ► time vs. space
- ► #examples vs. #features vs. #labels
- training vs. validation vs. prediction

• Size of the problem:

- Size of the problem:
 - # examples: $n > 10^6$

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - $\blacktriangleright~\#$ labels: $m>10^5$

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - # labels: $m > 10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{y} = \mathbf{W}^{\top} x$$

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - # labels: $m > 10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{y} = \mathbf{W}^{ op} x$$

► Train time complexity:

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - $\blacktriangleright~\#$ labels: $m>10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{m{y}} = \mathbf{W}^{ op} m{x}$$

• Train time complexity: $> 10^{17}$
- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - $\blacktriangleright~\#$ labels: $m>10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{\boldsymbol{y}} = \mathbf{W}^{ op} \boldsymbol{x}$$

- Train time complexity: $> 10^{17}$
- Space complexity:

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - $\blacktriangleright~\#$ labels: $m>10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{m{y}} = \mathbf{W}^{ op} m{x}$$

- Train time complexity: $> 10^{17}$
- Space complexity: $> 10^{11}$

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - $\blacktriangleright~\#$ labels: $m>10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{m{y}} = \mathbf{W}^{ op} m{x}$$

- Train time complexity: $> 10^{17}$
- ► Space complexity: > 10¹¹
- Test time complexity:

- Size of the problem:
 - # examples: $n > 10^6$
 - # features: $d > 10^6$
 - # labels: $m > 10^5$
- Naive solution: A dense linear model for each label (1-vs-All):

$$\hat{\boldsymbol{y}} = \mathbf{W}^{ op} \boldsymbol{x}$$

- Train time complexity: $> 10^{17}$
- ► Space complexity: > 10¹¹
- ▶ Test time complexity: > 10¹¹

• It does not have to be so hard:

- It does not have to be so hard:
 - ► High performance computing resources available

• It does not have to be so hard:

- High performance computing resources available
- Large data \longrightarrow sparse data (sparse features and labels)

• It does not have to be so hard:

- ► High performance computing resources available
- ► Large data → sparse data (sparse features and labels)
- ► Fast learning algorithms for standard learning problems exist

| Terminal | | | | 8 B) | \$ 🖶 40) sasan | 1 Mictoryfent CF | | | |
|--|----------|--------|----------|---------|----------------|------------------|------------------------|--|---|
| 0.912227 | 0.905463 | 22 | 22.0 | 1.0000 | -0.1043 | 87 | | | |
| 0.861865 | 0.811503 | 44 | 44.0 | -1.0000 | -0.0604 | 65 | | | |
| 0.823944 | 0.785142 | 87 | 87.0 | 1.0000 | -0.2309 | 60 | | | |
| 0.766675 | 0.709405 | 174 | 174.0 | 1.0000 | 0.0754 | 25 | | | |
| 0.642809 | 0.518943 | 348 | 348.0 | 1.0000 | 0.3440 | 47 | | | |
| 0.540082 | 0.437356 | 696 | 696.0 | 1.0000 | 0.9767 | 24 | | | |
| 0.450636 | 0.361190 | 1392 | 1392.0 | 1.0000 | 0.6204 | 181 | | | |
| 0.376935 | 0.303234 | 2784 | 2784.0 | 1.0000 | 0.4380 | 50 | | | |
| 0.320936 | 0.264938 | 5568 | 5568.0 | -1.0000 | -0.9257 | 89 | E 406 0 10000 0 4190 | | |
| 0.281048 | 0.241153 | 11135 | 11135.0 | 1.0000 | 1.0000 | 62 | 2004.1 0000 0001 | | |
| 0.249233 | 0.217415 | 22269 | 22269.0 | 1.0000 | 1.0000 | 140 | 11115.0 1.0000 UNC | | |
| 0.221765 | 0.194296 | 44537 | 44537.0 | 1.0000 | 1.0000 | 41 | 10071 1 1000 100 | | |
| 0.201490 | 0.181213 | 89073 | 89073.0 | -1.0000 | -1.0000 | 27 | 326231.0 1.0000 1.0000 | | |
| 0.187823 | 0.174157 | 178146 | 178146.0 | 1.0000 | 1.0000 | 49 | | | |
| 0.176267 | 0.164711 | 356291 | 356291.0 | -1.0000 | -1.0000 | 100 | | and the second s | |
| 0.165728 | 0.155188 | 712582 | 712582.0 | -1.0000 | -1.0000 | 69 | | | |
| finished run number of examples = 781265 weighted example sum = 7.8136+05 weighted label sum = -4.018e+04 average loss = 0.1645 best constant = -0.05143 total feature number = 59936409 vw -c rcv1.train.txt 1.46s user 0.21s system 189% cpu 0.883 total S:29PM 1-of-3-8: -/rcv1/norm [jl/ttypts/18] | | | | | | | | | |
| | | | | | | | | - 08:24 39:48 ◀ | - |

Vowpal Wabbit¹ at a lecture of John Langford²

¹ Vowpal Wabbit, http://hunch.net/~vw

² http://cilvr.cs.nyu.edu/doku.php?id=courses:bigdata:slides:start

Fast binary classification

- Data set: RCV1
- Predicted category: CCAT
- # training examples: 781 265
- # features: 60M
- Size: 1.1 GB
- Command line: time vw -sgd rcv1.train.txt -c
- Learning time: 1-3 secs on a laptop

• Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6~\text{secs}\approx11.6\text{-}34.7~\text{days}$

- Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6~\text{secs}\approx11.6\text{-}34.7~\text{days}$
- Can we further reduce computational costs?

- Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6~\text{secs}\approx11.6\text{-}34.7~\text{days}$
- Can we further reduce computational costs?
 - Smart 1-vs-All approaches

- Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6~\text{secs}\approx11.6\text{-}34.7~\text{days}$
- Can we further reduce computational costs?
 - Smart 1-vs-All approaches
 - Embeddings

- Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6~\text{secs}\approx11.6\text{-}34.7~\text{days}$
- Can we further reduce computational costs?
 - Smart 1-vs-All approaches
 - Embeddings
 - Tree-based methods

- Running Vowpal Wabbit for 10^6 labels will still require: $1\times10^6\text{-}3\times10^6\ \text{secs}\approx11.6\text{-}34.7\ \text{days}$
- Can we further reduce computational costs?
 - Smart 1-vs-All approaches
 - Embeddings
 - Tree-based methods
 - Label filtering/maximum inner product search (MIPS)

• Learning theory for an extremely large number of labels:

- Learning theory for an extremely large number of labels:
 - Statistical guarantees for the error rate that do not depend, or depend very weakly (sublinearly), on the total number of labels.

- Learning theory for an extremely large number of labels:
 - Statistical guarantees for the error rate that do not depend, or depend very weakly (sublinearly), on the total number of labels.
 - ► The **bound** on the error rate could be expressed in terms of the average number of **positive labels** (which is certainly much less than the total number of labels).

- Learning theory for an extremely large number of labels:
 - Statistical guarantees for the error rate that do not depend, or depend very weakly (sublinearly), on the total number of labels.
 - ► The bound on the error rate could be expressed in terms of the average number of positive labels (which is certainly much less than the total number of labels).
 - Particular performance guarantees depend on the considered loss function.

- Learning theory for an extremely large number of labels:
 - Statistical guarantees for the error rate that do not depend, or depend very weakly (sublinearly), on the total number of labels.
 - ► The bound on the error rate could be expressed in terms of the average number of positive labels (which is certainly much less than the total number of labels).
 - Particular performance guarantees depend on the considered loss function.
 - Different theoretical settings: statistical learning theory, learning reductions, online learning.

• Training and prediction under limited time and space budget:

- Training and prediction under limited time and space budget:
 - Restricted computational resources (time and space) for both training and prediction.

- Training and prediction under limited time and space budget:
 - Restricted computational resources (time and space) for both training and prediction.
 - ► A trade-off between computational (time and space) complexity and the predictive performance.

- Training and prediction under limited time and space budget:
 - Restricted computational resources (time and space) for both training and prediction.
 - ► A trade-off between computational (time and space) complexity and the predictive performance.
 - By imposing hard constraints on time and space budget, the challenge is then to optimize the predictive performance of an algorithm under these constraints.

• Unreliable learning information:

- Unreliable learning information:
 - We cannot expect that all labels will be properly checked and assigned to training examples.

- Unreliable learning information:
 - We cannot expect that all labels will be properly checked and assigned to training examples.
 - Therefore we often deal with a problem of learning with missing labels or learning from positive and unlabeled examples.

• Performance measures:

- Performance measures:
 - ► Typical performance measures such as 0/1 or Hamming loss do not fit well to the extreme setting.

- Performance measures:
 - ➤ Typical performance measures such as 0/1 or Hamming loss do not fit well to the extreme setting.
 - Other measures are often used such as **precision@k** or the **F**-measure.

- Performance measures:
 - ➤ Typical performance measures such as 0/1 or Hamming loss do not fit well to the extreme setting.
 - Other measures are often used such as precision@k or the F-measure.
 - ► However, it remains an **open question** how to **design loss functions** suitable for extreme classification.

Do we search in the right place?



Figure: ³ A similar comics has been earlier used by Asela Gunawardana.⁴

³ Source: Florence Morning News, Mutt and Jeff Comic Strip, Page 7, Florence, South Carolina,1942

⁴ Asela Gunawardana, Evaluating Machine Learned User Experiences. Extreme Classification Workshop. NIPS 2015

• Long-tail label distributions and zero-shot learning:

- Long-tail label distributions and zero-shot learning:
 - A close relation to the problem of estimating distributions over large alphabets.

- Long-tail label distributions and zero-shot learning:
 - A close relation to the problem of estimating distributions over large alphabets.
 - ► The distribution of label frequencies is often characterized by a **long-tail** for which proper **smoothing** (like add-constant or Good-Turing estimates) or **calibration** techniques (like isotonic regression or domain adaptation) have to be used.
Statistical challenges

- Long-tail label distributions and zero-shot learning:
 - A close relation to the problem of estimating distributions over large alphabets.
 - ► The distribution of label frequencies is often characterized by a **long-tail** for which proper **smoothing** (like add-constant or Good-Turing estimates) or **calibration** techniques (like isotonic regression or domain adaptation) have to be used.
 - ► In practical applications, learning algorithms run in rapidly changing environments: new labels may appear during testing/prediction phase (⇒ zero-shot learning)

Statistical challenges

- Long-tail label distributions and zero-shot learning:
 - ► Frequency of labels in the WikiLSHTC dataset:⁵



• Many labels with only few examples (\Rightarrow one-shot learning).

⁵ http://manikvarma.org/downloads/XC/XMLRepository.html

Statistical challenges

- Long-tail label distributions and zero-shot learning:
 - Frequency of frequencies for the WikiLSHTC dataset:



► The missing mass obtained by the Good-Turing estimate: 0.014.

Object of interest: documents, images, ads, posts ...











27 / 94

Test example x







True outcome y





Training data $\{oldsymbol{x}_i,oldsymbol{y}_i\}_1^n$















- Input $x \in \mathcal{X}$ drawn from a distribution $\mathbf{P}(x)$.
 - usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.

- Input $x \in \mathcal{X}$ drawn from a distribution $\mathbf{P}(x)$.
 - usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- Outcome $y \in \mathcal{Y}$ drawn from a distribution $\mathbf{P}(y \,|\, x)$.
 - a vector of labels $\boldsymbol{y} = (y_1, y_2, \dots, y_m)$.

- Input $x \in \mathcal{X}$ drawn from a distribution $\mathbf{P}(x)$.
 - usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- Outcome $y \in \mathcal{Y}$ drawn from a distribution $\mathbf{P}(y \,|\, x)$.
 - a vector of labels $\boldsymbol{y} = (y_1, y_2, \dots, y_m)$.
- Prediction $\hat{y} = h(x)$ by means of prediction function $h: \mathcal{X} \to \mathcal{Y}$.
 - h returns prediction $\hat{y} = h(x)$ for every input x.

- Input $x \in \mathcal{X}$ drawn from a distribution $\mathbf{P}(x)$.
 - usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- Outcome $y \in \mathcal{Y}$ drawn from a distribution $\mathbf{P}(y \,|\, x)$.
 - a vector of labels $\boldsymbol{y} = (y_1, y_2, \dots, y_m)$.
- Prediction $\hat{y} = h(x)$ by means of prediction function $h: \mathcal{X} \to \mathcal{Y}$.
 - h returns prediction $\hat{y} = h(x)$ for every input x.
- Loss of our prediction: $\ell(\boldsymbol{y}, \hat{\boldsymbol{y}})$.
 - $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ is a task-specific loss function.

- Input $x \in \mathcal{X}$ drawn from a distribution $\mathbf{P}(x)$.
 - usually a feature vector, $\mathcal{X} \subseteq \mathbb{R}^d$.
- Outcome $y \in \mathcal{Y}$ drawn from a distribution $\mathbf{P}(y \,|\, x)$.
 - a vector of labels $\boldsymbol{y} = (y_1, y_2, \dots, y_m)$.
- Prediction $\hat{y} = h(x)$ by means of prediction function $h: \mathcal{X} \to \mathcal{Y}$.
 - h returns prediction $\hat{y} = h(x)$ for every input x.
- Loss of our prediction: $\ell({m y}, \hat{{m y}}).$
 - $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$ is a task-specific loss function.
- Goal: find a prediction function with small loss.

• Goal: minimize the expected loss over all examples (risk):

$$L_{\ell}(h) = \mathbb{E}_{(\boldsymbol{x}, \boldsymbol{y}) \sim \mathbf{P}} \left[\ell(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) \right].$$

• Goal: minimize the expected loss over all examples (risk):

$$L_{\ell}(h) = \mathbb{E}_{(\boldsymbol{x}, \boldsymbol{y}) \sim \mathbf{P}} \left[\ell(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) \right].$$

• The **optimal** prediction function over all possible functions expressed conditionally for a given *x*:

$$h^*(x) = \operatorname*{arg\,min}_h L_\ell(h|x),$$

(so called **Bayes prediction function**).

• Hamming loss:

$$\ell_H(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) = \frac{1}{m} \sum_{j=1}^m \llbracket y_j \neq h_j(\boldsymbol{x}) \rrbracket,$$

⁶ K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. On loss minimization and label dependence in multi-label classification. *Machine Learning*, 88:5–45, 2012 31/94

• Hamming loss:

$$\ell_H(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) = \frac{1}{m} \sum_{j=1}^m \llbracket y_j \neq h_j(\boldsymbol{x}) \rrbracket,$$

• Sparse labels \Rightarrow Hamming loss of an all-zero classifier close to **0**.

⁶ K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. On loss minimization and label dependence in multi-label classification. *Machine Learning*, 88:5–45, 2012 31 / 94

• Hamming loss:

$$\ell_H(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) = \frac{1}{m} \sum_{j=1}^m \llbracket y_j \neq h_j(\boldsymbol{x}) \rrbracket,$$

Sparse labels ⇒ Hamming loss of an all-zero classifier close to 0.
 The optimal strategy:⁶

$$h_j^*(\pmb{x}) = [\![\eta_j(\pmb{x})>0.5]\!]\,,$$
 where $\eta_j(\pmb{x}) = \mathbf{P}(y_j=1\,|\,\pmb{x}).$

⁶ K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. On loss minimization and label dependence in multi-label classification. *Machine Learning*, 88:5–45, 2012 31/94

• Hamming loss:

where η_i

$$\ell_H(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) = \frac{1}{m} \sum_{j=1}^m \llbracket y_j \neq h_j(\boldsymbol{x}) \rrbracket,$$

Sparse labels ⇒ Hamming loss of an all-zero classifier close to 0.
 The optimal strategy:⁶

$$h_j^*(m{x}) = \llbracket \eta_j(m{x}) > 0.5
rbracket \,,$$

 $(m{x}) = \mathbf{P}(y_j = 1 \mid m{x}).$
 $\hat{\eta}_1(m{x}) \ \hat{\eta}_2(m{x}) \ \hat{\eta}_3(m{x}) \ \hat{\eta}_4(m{x}) \ \hat{\eta}_5(m{x}) \ \hat{\eta}_6(m{x}) \ \hat{\eta}_7(m{x})$
 $au = 0.5$

⁶ K. Dembczyński, W. Waegeman, W. Cheng, and E. Hüllermeier. On loss minimization and label dependence in multi-label classification. *Machine Learning*, 88:5–45, 2012 31/94
Precision

• Precision at position k:

$$\operatorname{prec}@k(\boldsymbol{y}, \boldsymbol{h}, \boldsymbol{x}) = rac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \llbracket y_j = 1
rbracket,$$

where $\hat{\mathcal{Y}}_k$ is a set of k labels predicted by h.

Precision

• Precision at position k:

$$\operatorname{prec}@k(\boldsymbol{y}, \boldsymbol{h}, \boldsymbol{x}) = \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \llbracket y_j = 1 \rrbracket,$$

where $\hat{\mathcal{Y}}_k$ is a set of k labels predicted by h.

• The optimal strategy: select top k labels according to $\eta_j(x)$.

Precision

• **Precision** at position k:

$$\operatorname{prec}@k(\boldsymbol{y}, \boldsymbol{h}, \boldsymbol{x}) = \frac{1}{k} \sum_{j \in \hat{\mathcal{Y}}_k} \llbracket y_j = 1 \rrbracket,$$

where $\hat{\mathcal{Y}}_k$ is a set of k labels predicted by h.

• The optimal strategy: select top k labels according to $\eta_j(x)$.



Normalized Discounted Cumulative Gain

• Normalized Discounted Cumulative Gain at position k:

NDCG@
$$k(\boldsymbol{y}, f, \boldsymbol{x}) = N_k(\boldsymbol{y}) \sum_{r=1}^k \frac{y_{\sigma(r)}}{\log(1+r)},$$

where σ is a permutation of labels for x returned by ranker f, and $N_k(y)$ normalizes NDCG@k to the interval [0, 1]:

$$N_k(\boldsymbol{y}) = \left(\sum_{r=1}^{\max(k,\sum_{i=1}^m y_i)} \frac{1}{\log(1+r)}\right)^{-1}$$

Normalized Discounted Cumulative Gain

• **The optimal strategy**: rank labels according to the following marginal quantities:

$$\Delta_j(\boldsymbol{x}) = \sum_{\boldsymbol{y}: y_j = 1} N_k(\boldsymbol{y}) \mathbf{P}(\boldsymbol{y} \,|\, \boldsymbol{x})$$

Normalized Discounted Cumulative Gain

• The optimal strategy: rank labels according to the following marginal quantities:

$$\Delta_j(oldsymbol{x}) = \sum_{oldsymbol{y}: y_j = 1} N_k(oldsymbol{y}) \mathbf{P}(oldsymbol{y} \mid oldsymbol{x})$$
 $\Delta_6(oldsymbol{x}) \; \Delta_2(oldsymbol{x}) \; \Delta_5(oldsymbol{x}) \; \Delta_1(oldsymbol{x}) \; \Delta_7(oldsymbol{x}) \; \Delta_3(oldsymbol{x}) \; \Delta_4(oldsymbol{x})$

• The macro F-measure (F-score):

 y_{11}

 y_{21}

 y_{31}

 y_{41}

 y_{51}

 y_{61}

$$F_M(\boldsymbol{Y}, \widehat{\boldsymbol{Y}}) = \frac{1}{m} \sum_{j=1}^m F(\boldsymbol{y}_{\cdot j}, \widehat{\boldsymbol{y}}_{\cdot j}) = \frac{1}{m} \sum_{j=1}^m \frac{2\sum_{i=1}^n y_{ij} \widehat{y}_{ij}}{\sum_{i=1}^n y_{ij} + \sum_{i=1}^n \widehat{y}_{ij}}$$

True labels

 y_{13}

 y_{23}

 y_{33}

 y_{43}

 y_{53}

 y_{63}

 y_{14}

 y_{24}

 y_{34}

 y_{44}

 y_{54}

 y_{64}

 y_{12}

 y_{22}

 y_{32}

 y_{42}

 y_{52}

 y_{62}

| Pred | licted | labels |
|------|--------|--------|
| | | |

| \hat{y}_{11} | \hat{y}_{12} | \hat{y}_{13} | \hat{y}_{14} |
|----------------|----------------|----------------|----------------|
| \hat{y}_{21} | \hat{y}_{22} | \hat{y}_{23} | \hat{y}_{24} |
| \hat{y}_{31} | \hat{y}_{32} | \hat{y}_{33} | \hat{y}_{34} |
| \hat{y}_{41} | \hat{y}_{42} | \hat{y}_{43} | \hat{y}_{44} |
| \hat{y}_{51} | \hat{y}_{52} | \hat{y}_{53} | \hat{y}_{54} |
| \hat{y}_{61} | \hat{y}_{62} | \hat{y}_{63} | \hat{y}_{64} |

• The macro F-measure (F-score):

$$F_M(\boldsymbol{Y}, \widehat{\boldsymbol{Y}}) = \frac{1}{m} \sum_{j=1}^m F(\boldsymbol{y}_{\cdot j}, \widehat{\boldsymbol{y}}_{\cdot j}) = \frac{1}{m} \sum_{j=1}^m \frac{2\sum_{i=1}^n y_{ij} \widehat{y}_{ij}}{\sum_{i=1}^n y_{ij} + \sum_{i=1}^n \widehat{y}_{ij}}$$

True labels

| y_{11} | y_{12} | y_{13} | y_{14} |
|----------|----------|----------|----------|
| y_{21} | y_{22} | y_{23} | y_{24} |
| y_{31} | y_{32} | y_{33} | y_{34} |
| y_{41} | y_{42} | y_{43} | y_{44} |
| y_{51} | y_{52} | y_{53} | y_{54} |
| y_{61} | y_{62} | y_{63} | y_{64} |

| \hat{y}_{11} | \hat{y}_{12} | \hat{y}_{13} | \hat{y}_{14} |
|----------------|----------------|----------------|----------------|
| \hat{y}_{21} | \hat{y}_{22} | \hat{y}_{23} | \hat{y}_{24} |
| \hat{y}_{31} | \hat{y}_{32} | \hat{y}_{33} | \hat{y}_{34} |
| \hat{y}_{41} | \hat{y}_{42} | \hat{y}_{43} | \hat{y}_{44} |
| \hat{y}_{51} | \hat{y}_{52} | \hat{y}_{53} | \hat{y}_{54} |
| \hat{y}_{61} | \hat{y}_{62} | \hat{y}_{63} | \hat{y}_{64} |

• The macro F-measure (F-score):

 y_{11}

 y_{21}

 y_{31}

 y_{41}

 y_{51}

 y_{61}

$$F_M(\boldsymbol{Y}, \widehat{\boldsymbol{Y}}) = \frac{1}{m} \sum_{j=1}^m F(\boldsymbol{y}_{\cdot j}, \widehat{\boldsymbol{y}}_{\cdot j}) = \frac{1}{m} \sum_{j=1}^m \frac{2\sum_{i=1}^n y_{ij} \widehat{y}_{ij}}{\sum_{i=1}^n y_{ij} + \sum_{i=1}^n \widehat{y}_{ij}}$$

True labels

 y_{13}

 y_{23}

 y_{33}

 y_{43}

 y_{53}

 y_{63}

 y_{14}

 y_{24}

 y_{34}

 y_{44}

 y_{54}

 y_{64}

 y_{12}

 y_{22}

 y_{32}

 y_{42}

 y_{52}

 y_{62}

| Pred | licted | labels |
|------|--------|--------|
| | | |

| \hat{y}_{11} | \hat{y}_{12} | \hat{y}_{13} | \hat{y}_{14} |
|----------------|----------------|----------------|----------------|
| \hat{y}_{21} | \hat{y}_{22} | \hat{y}_{23} | \hat{y}_{24} |
| \hat{y}_{31} | \hat{y}_{32} | \hat{y}_{33} | \hat{y}_{34} |
| \hat{y}_{41} | \hat{y}_{42} | \hat{y}_{43} | \hat{y}_{44} |
| \hat{y}_{51} | \hat{y}_{52} | \hat{y}_{53} | \hat{y}_{54} |
| \hat{y}_{61} | \hat{y}_{62} | \hat{y}_{63} | \hat{y}_{64} |

• The macro F-measure (F-score):

 y_{11}

 y_{21}

 y_{31}

 y_{41}

 y_{51}

 y_{61}

$$F_M(\boldsymbol{Y}, \widehat{\boldsymbol{Y}}) = \frac{1}{m} \sum_{j=1}^m F(\boldsymbol{y}_{\cdot j}, \widehat{\boldsymbol{y}}_{\cdot j}) = \frac{1}{m} \sum_{j=1}^m \frac{2\sum_{i=1}^n y_{ij} \widehat{y}_{ij}}{\sum_{i=1}^n y_{ij} + \sum_{i=1}^n \widehat{y}_{ij}}$$

True labels

 y_{13}

 y_{23}

 y_{33}

 y_{43}

 y_{53}

 y_{63}

 y_{14}

 y_{24}

 y_{34}

 y_{44}

 y_{54}

 y_{64}

 y_{12}

 y_{22}

 y_{32}

 y_{42}

 y_{52}

 y_{62}

| Predicted label |
|-----------------|
|-----------------|

| \hat{y}_{11} | \hat{y}_{12} | \hat{y}_{13} | \hat{y}_{14} |
|----------------|----------------|----------------|----------------|
| \hat{y}_{21} | \hat{y}_{22} | \hat{y}_{23} | \hat{y}_{24} |
| \hat{y}_{31} | \hat{y}_{32} | \hat{y}_{33} | \hat{y}_{34} |
| \hat{y}_{41} | \hat{y}_{42} | \hat{y}_{43} | \hat{y}_{44} |
| \hat{y}_{51} | \hat{y}_{52} | \hat{y}_{53} | \hat{y}_{54} |
| \hat{y}_{61} | \hat{y}_{62} | \hat{y}_{63} | \hat{y}_{64} |

• The macro F-measure (F-score):

 y_{11}

 y_{21}

 y_{31}

 y_{41}

 y_{51}

 y_{61}

$$F_M(\boldsymbol{Y}, \widehat{\boldsymbol{Y}}) = \frac{1}{m} \sum_{j=1}^m F(\boldsymbol{y}_{\cdot j}, \widehat{\boldsymbol{y}}_{\cdot j}) = \frac{1}{m} \sum_{j=1}^m \frac{2\sum_{i=1}^n y_{ij} \widehat{y}_{ij}}{\sum_{i=1}^n y_{ij} + \sum_{i=1}^n \widehat{y}_{ij}}$$

True labels

 y_{13}

 y_{23}

 y_{33}

 y_{43}

 y_{53}

 y_{63}

 y_{14}

 y_{24}

 y_{34}

 y_{44}

 y_{54}

 y_{64}

 y_{12}

 y_{22}

 y_{32}

 y_{42}

 y_{52}

 y_{62}

| Pred | licted | labels |
|------|--------|--------|
| | | |

| \hat{y}_{11} | \hat{y}_{12} | \hat{y}_{13} | \hat{y}_{14} |
|----------------|----------------|----------------|----------------|
| \hat{y}_{21} | \hat{y}_{22} | \hat{y}_{23} | \hat{y}_{24} |
| \hat{y}_{31} | \hat{y}_{32} | \hat{y}_{33} | \hat{y}_{34} |
| \hat{y}_{41} | \hat{y}_{42} | \hat{y}_{43} | \hat{y}_{44} |
| \hat{y}_{51} | \hat{y}_{52} | \hat{y}_{53} | \hat{y}_{54} |
| \hat{y}_{61} | \hat{y}_{62} | \hat{y}_{63} | \hat{y}_{64} |

• Can be solved by reduction to m independent binary problems.⁷

⁷ O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon. Consistent multilabel classification. In NIPS, 2015

- Can be solved by reduction to *m* independent binary problems.⁷
- Thresholding the conditional probabilities:

$$F(\tau) = \frac{2\int_{\mathcal{X}} \eta(\boldsymbol{x}) \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}{\int_{\mathcal{X}} \eta(\boldsymbol{x}) \, \mathrm{d}\mu(\boldsymbol{x}) + \int_{\mathcal{X}} \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}$$

⁷ O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon. Consistent multilabel classification. In NIPS, 2015

- Can be solved by reduction to m independent binary problems.⁷
- Thresholding the conditional probabilities:

$$F(\tau) = \frac{2\int_{\mathcal{X}} \eta(\boldsymbol{x}) \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}{\int_{\mathcal{X}} \eta(\boldsymbol{x}) \, \mathrm{d}\mu(\boldsymbol{x}) + \int_{\mathcal{X}} \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}.$$

• The optimal F-measure is $F(\tau^*)$: no binary classifier can be better.

⁷ O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon. Consistent multilabel classification. In NIPS, 2015

- Can be solved by reduction to m independent binary problems.⁷
- Thresholding the conditional probabilities:

$$F(\tau) = \frac{2\int_{\mathcal{X}} \eta(\boldsymbol{x}) \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}{\int_{\mathcal{X}} \eta(\boldsymbol{x}) \, \mathrm{d}\mu(\boldsymbol{x}) + \int_{\mathcal{X}} \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}.$$

- The optimal F-measure is $F(\tau^*)$: no binary classifier can be better.
- The optimal solution satisfies the following condition: $F(\tau^*) = 2\tau^*$.

⁷ O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon. Consistent multilabel classification. In NIPS, 2015

- Can be solved by reduction to m independent binary problems.⁷
- Thresholding the conditional probabilities:

$$F(\tau) = \frac{2 \int_{\mathcal{X}} \eta(\boldsymbol{x}) \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}{\int_{\mathcal{X}} \eta(\boldsymbol{x}) \, \mathrm{d}\mu(\boldsymbol{x}) + \int_{\mathcal{X}} \llbracket \eta(\boldsymbol{x}) \ge \tau \rrbracket \, \mathrm{d}\mu(\boldsymbol{x})}.$$

- The optimal F-measure is $F(\tau^*)$: no binary classifier can be better.
- The optimal solution satisfies the following condition: $F(\tau^*) = 2\tau^*$.

$$\hat{\eta}_1(x) \quad \hat{\eta}_2(x) \quad \hat{\eta}_3(x) \quad \hat{\eta}_4(x) \quad \hat{\eta}_5(x) \quad \hat{\eta}_6(x) \quad \hat{\eta}_7(x)$$



⁷ O. Koyejo, N. Natarajan, P. Ravikumar, and I. Dhillon. Consistent multilabel classification. In NIPS, 2015

Predictive model

• From the above analysis we can conclude:

Predictive model

• From the above analysis we can conclude:

We need to train models that accurately estimate marginal probabilities or other related marginal quantities

Agenda

- 1 Extreme classification: applications and challenges
- 2 Algorithms
- 3 Live demonstration

Label embedding methods



- Shallow Networks SVM
 - Direct mapping of input to output
- Output Input layer layer \boldsymbol{y}_1 $\boldsymbol{x}^{(1)}$ y_2 $oldsymbol{x}^{(2)}$ $\boldsymbol{x}^{(3)}$ y_4 $x^{(100\,000)}$ ${m y}_{670,000}$

- Label embedding
 - Mapping input to output via embedding layer



Let ${\boldsymbol Y}$ be the label matrix of dimensionality $m\times n$ such that each row denotes the labels in an instance

• Label embedding methods - Assume a low rank structure in the label matrix Y i.e. the m columns of the Y can be effectively represented by $\hat{m} \ll m$ columns

Let ${\boldsymbol Y}$ be the label matrix of dimensionality $m\times n$ such that each row denotes the labels in an instance

• Label embedding methods - Assume a low rank structure in the label matrix Y i.e. the m columns of the Y can be effectively represented by $\hat{m} \ll m$ columns



Let ${\boldsymbol Y}$ be the label matrix of dimensionality $m\times n$ such that each row denotes the labels in an instance

• Label embedding methods - Assume a low rank structure in the label matrix Y i.e. the m columns of the Y can be effectively represented by $\hat{m} \ll m$ columns



• Each label vector $\boldsymbol{y}_i \in \{0,1\}^m$ is hence projected to $\boldsymbol{z}_i \in \mathbb{R}^{\hat{m}}$ by a projection matrix $\boldsymbol{U} \in \mathbb{R}^{m \times \hat{m}}$, such that $\boldsymbol{U} \boldsymbol{y}_i = \boldsymbol{z}_i$

Let ${\boldsymbol Y}$ be the label matrix of dimensionality $m\times n$ such that each row denotes the labels in an instance

• Label embedding methods - Assume a low rank structure in the label matrix Y i.e. the m columns of the Y can be effectively represented by $\hat{m} \ll m$ columns



- Each label vector $\boldsymbol{y}_i \in \{0,1\}^m$ is hence projected to $\boldsymbol{z}_i \in \mathbb{R}^{\hat{m}}$ by a projection matrix $\boldsymbol{U} \in \mathbb{R}^{m \times \hat{m}}$, such that $\boldsymbol{U} \boldsymbol{y}_i = \boldsymbol{z}_i$
- Regressors $oldsymbol{V}$ are then trained to predict $oldsymbol{z}_i = oldsymbol{V} oldsymbol{x}_i$

- Otherwise, we may directly map the input $m{x}$ to its corresponding prediction $\hat{m{y}} = m{U}^\dagger m{V} m{x}$
- Here U^{\dagger} is the decompression matrix which lifts intermediary prediction Vx from the embedding space to original label space

- Otherwise, we may directly map the input x to its corresponding prediction $\hat{y} = U^{\dagger}Vx$
- Here U^{\dagger} is the decompression matrix which lifts intermediary prediction Vx from the embedding space to original label space



Label Embedding Methods

Various methods differ in terms of the embedding/compression and de-compression methods they employ

- Methods for label embedding
 - Compressed Sensing ⁸
 - Label Embedding for Missing Labels (LEML) ⁹
 - Sparse Local Embeddings for Extreme Classification (SLEEC) ¹⁰

D. Hsu, S. Kakade, J. Langford, and T. Zhang. Multi-label prediction via compressed sensing. In NIPS, 2009

⁹ Hsiang-Fu Yu, Prateek Jain, Purushottam Kar, and Inderjit S. Dhillon. Large-scale Multi-label Learning with Missing Labels. In ICML, 2014

¹⁰ Kush Bhatia, Himanshu Jain, Purushottam Kar, Manik Varma, and Prateek Jain. Sparse local embeddings for extreme multilabel classification. In NIPS, 2015

LEML learns global embeddings under the assumption :

• The label matrix is low rank

LEML learns global embeddings under the assumption :

- The label matrix is low rank
- Preserve global distances

LEML learns global embeddings under the assumption :

- The label matrix is low rank
- Preserve global distances



LEML Optimization - Direct mapping of input to output





LEML Optimization - Direct mapping of input to output



• The above optimization problem is non-convex

LEML Optimization - Direct mapping of input to output



- The above optimization problem is non-convex
- Fixing one of U^{\dagger} or V leads convex optimization

LEML Optimization - Direct mapping of input to output



- The above optimization problem is non-convex
- Fixing one of U^{\dagger} or V leads convex optimization
- Alternating minimization over \boldsymbol{U}^{\dagger} and \boldsymbol{V} is performed
LEML - Pros and Cons

- Advantages of label embedding
 - Exploit label correlations
 - Ease of implementation

LEML - Pros and Cons

- Advantages of label embedding
 - Exploit label correlations
 - Ease of implementation
- Disadvantages of label embedding
 - Assuption fails for long tail of labels
 - Prediction complexity $\Omega(\hat{m}(m+\hat{d}))$

• The embedding matrix $Z \in \mathbb{R}^{\hat{m} imes n}$ into which the original labels are projected is learnt by

$$\min_{\boldsymbol{Z} \in \hat{m} \times n} ||P_{\Omega}(\boldsymbol{Y}^{T}\boldsymbol{Y}) - P_{\Omega}(\boldsymbol{Z}^{T}\boldsymbol{Z})||_{F}^{2} + \lambda ||\boldsymbol{Z}||_{1}$$

where Ω denotes the index set of nearest neighbors $(i, j) \in \Omega$ iff $j \in \mathcal{N}_i$, and \mathcal{N}_i is the set of nearest neighbors of instance i.

• The embedding matrix $Z \in \mathbb{R}^{\hat{m} imes n}$ into which the original labels are projected is learnt by

$$\min_{\boldsymbol{Z} \in \hat{m} \times n} ||P_{\Omega}(\boldsymbol{Y}^{T}\boldsymbol{Y}) - P_{\Omega}(\boldsymbol{Z}^{T}\boldsymbol{Z})||_{F}^{2} + \lambda ||\boldsymbol{Z}||_{1}$$

where Ω denotes the index set of nearest neighbors $(i, j) \in \Omega$ iff $j \in \mathcal{N}_i$, and \mathcal{N}_i is the set of nearest neighbors of instance i.



• Sparsity in embedding is obtained by ℓ_1 regularization $||m{Z}||_1$

- Sparsity in embedding is obtained by ℓ_1 regularization $||m{Z}||_1$
- Locality in the embedding is based on following a nearest neighbors approach and is related to N_i given by

$$\mathcal{N}_i = rgmax_{S,|S| \le lpha.n} \sum_{j \in S} (\boldsymbol{Y}_i^T \boldsymbol{Y}_j)$$



Overall optimization

• The overall SLEEC objective is therefore given by

 $\min_{\boldsymbol{V} \in \mathbb{R}^{\hat{m} \times d}} ||P_{\Omega}(\boldsymbol{Y}^{T}\boldsymbol{Y}) - P_{\Omega}(\boldsymbol{X}^{T}\boldsymbol{V}^{T}\boldsymbol{V}\boldsymbol{X})||_{F}^{2} + \lambda ||\boldsymbol{V}||_{F}^{2} + \mu ||\boldsymbol{V}\boldsymbol{X}||_{1}$

• To optimize above is challenging due to non-convexity, non-differentiability and scale of the problem

• Simplification

 \blacktriangleright Finding the embedding matrix, without ℓ_1 regularization

$$\min_{\boldsymbol{Z} \in \hat{m} \times n} ||P_{\Omega}(\boldsymbol{Y}^{T}\boldsymbol{Y}) - P_{\Omega}(\boldsymbol{Z}^{T}\boldsymbol{Z})||_{F}^{2} \equiv \min_{\boldsymbol{M} \succeq 0, rank(\boldsymbol{M}) \leq \hat{m}} ||P_{\Omega}(\boldsymbol{Y}^{T}\boldsymbol{Y}) - P_{\Omega}(\boldsymbol{M})|$$

where $oldsymbol{M} = oldsymbol{Z}^T oldsymbol{Z}$

► After learning Z in the above step,

$$\min_{oldsymbol{V}\in\mathbb{R}^{\hat{m} imes d}}||oldsymbol{Z}-oldsymbol{V}oldsymbol{X}||_F^2+\lambda||oldsymbol{V}||_F^2+\mu||oldsymbol{V}oldsymbol{X}||_1$$

Smart 1-vs-all approaches

Learning Linear Decision Boundaries



Learning process typically involves finding the optimal linear separator as follows:

• Template of the optimization problem being solved

$$\hat{\boldsymbol{w}} = \arg\min_{\boldsymbol{w}} R_{emp}(\boldsymbol{w}) + \lambda \ Reg(\boldsymbol{w})$$

where Reg(w) is the regularization term to avoid complex models and $R_{emp}(.)$ represents the empirical

Learning Linear Decision Boundaries



Learning process typically involves finding the optimal linear separator as follows:

• For Support Vector Machine - using squared hinge loss

$$\arg\min_{\boldsymbol{w}} \frac{\lambda}{2} ||\boldsymbol{w}||^2 + \sum_{i=1}^n (\max(0, 1 - y_i \boldsymbol{w}^T \boldsymbol{x}_i))^2$$

Bag of Words Data Representation

Top-3 training instances and their corresponding labels from a Wikipedia Subset

1,2,3,4,5 1:0.1498 2:0.1216 3:0.0521 4:0.0637 5:0.0517 6:0.0891 7:0.0994 8:0.1944 9:0.3821 10:0.093 8 11:0.0361 12:0.186 13:0.0859 14:0.0641 15:0.0616 16:0.0821 17:0.1085 18:0.0636 19:0.0714 20:0.076 8 21:0.1532 22:0.0491 23:0.1017 24:0.1542 25:0.1136 26:0.0965 27:0.0473 28:0.1223 29:0.0936 30:0.09 29 31:0.1665 32:0.1444 33:0.0581 34:0.0622 35:0.0578 36:0.1462 37:0.0704 38:0.0623 39:0.2466 40:0.1 539 41:0.0992 42:0.071 43:0.1106 44:0.0849 45:0.1113 46:0.0725 47:0.1093 48:0.073 49:0.1415 50:0.09 23 51:0.0547 52:0.1022 53:0.1009 54:0.0431 55:0.0889 56:0.1315 57:0.0637 58:0.1108 59:0.0656 60:0.0 873 61:0.0813 62:0.0864 63:0.2023 64:0.2295 65:0.067 66:0.067 67:0.0614 68:0.0447 69:0.0878 70:0.06 71:0.0711 72:0.0755 73:0.0598 74:0.1457 75:0.058 76:0.0862 77:0.0748 78:0.0351 79:0.0552 477142:1

6,7,8,9 14:0.0562 27:0.0414 28:0.0633 42:0.0623 51:0.1006 54:0.0639 63:0.068 66:0.0587 71:0.0624 80 10.0799 81:0.0477 82:0.1125 83:0.0356 84:0.0627 85:0.0989 86:0.0517 87:0.0395 88:0.3179 89:0.1795 9 0:0.0576 91:0.0916 92:0.0349 93:0.069 94:0.1272 95:0.0683 96:0.0836 97:0.0427 98:0.0502 99:0.046 10 0:0.1292 101:0.0429 102:0.048 103:0.0469 104:0.0839 105:0.0737 106:0.1424 107:0.0534 108:0.0818 109 1:0.0804 110:0.0917 111:0.1455 112:0.1106 113:0.1168 114:0.0621 115:0.0995 116:0.0526 117:0.0891 118 1:0.0714 119:0.0529 120:0.0555 121:0.0891 122:0.07 123:0.0525 124:0.1049 125:0.0255 127:0.0451 127:0 0.789 128:0.0503 129:0.0959 130:0.3668 131:0.0821 132:0.0715 133:0.2504 134:0.0647 135:0.0521 136:0 0.988 137:0.0579 138:0.0449 139:0.0682 140:0.0683 141:0.0631 142:0.0797 143:0.0615 144:0.0826 145:0 0.436 155:0.0608 147:0.0629 148:0.0724 149:0.0884 150:0.0633 151:0.0537 152:0.0481 153:0.1023 154:0. 0436 155:0.0609 156:0.0338 157:0.0555 158:0.0628 159:0.1355 160:0.221 161:0.0468 162:0.0684 163:0.0 659 164:0.0483 165:0.0884 166:0.063 167:0.1083 168:0.0503 169:0.0853 170:0.0794 171:0.0364 172:0.07 8 132:0.0414 174:0.0744 175:0.1535 176:0.0482 177:0.0417 178:0.0629 179:0.083 180:0.0709 181:0.050 8 182:0.1009 477142:1

10,11,12,13,14,15,16 34:0.0992 125:0.0463 173:0.0752 183:0.0718 184:0.0924 185:0.219 186:0.0927 187 10.1781 188:0.1625 189:0.135 190:0.0868 191:0.1272 192:0.103 193:0.1285 194:0.1504 195:0.1447 196:0 .1869 197:0.3767 198:0.1646 199:0.4598 200:0.2042 201:0.0736 202:0.1852 203:0.1223 204:0.0782 205:0 .2479 206:0.0893 207:0.1112 208:0.0861 209:0.1697 210:0.1213 211:0.2672 212:0.1393 477142:1

Computational Challenge - Big Data

| | #Labels (m) | #Features (d) | #Instances (n) | Parameters $(m \times d)$ | Training Data | Model size |
|-------------------------|---------------|---------------|------------------|---------------------------|---------------|------------|
| DMOZ ¹¹ | 12,294 | 347,256 | 93,805 | 4,269,165,264 | 125MB | 17 GB |
| Wikipedia | 325,056 | 1,617,899 | 1,778,351 | 525,907,777,344 | 1GB | 870GB |
| Delicious ¹² | 205,443 | 782,585 | 196,606 | 160,776,610,155 | 1GB | 350GB |
| Amazon | 670,091 | 135,909 | 490,449 | 91,071,397,719 | 600MB | 250GB |

- Naive application of binary one-vs-rest linear classification
 - Billions of parameters Computational complexity
 - ► TeraBytes of disk space to store the model Space Complexity

¹¹ From LSHTC challenge

 $^{^{12}\,\}mathrm{From}$ XMC repository

Distribution of Learnt Weights

Weights $(W_{d',m'})$ learnt by using One-vs-rest SVM for Wiki-31K dataset (101k dimensional data with 31k labels) from XMC repository

$$\min_{\boldsymbol{w}_{m'}} \left[||\boldsymbol{w}_{m'}||_2^2 + C \sum_{i=1}^n (\max(0, 1 - s_{m'_i} \boldsymbol{w}_{m'}^T \boldsymbol{x}_i))^2 \right]$$

 $s_{m'_i} = +1$ if instance *i* has label m'_i , -1 otherwise .





Figure: Distribution of learnt weights before pruning

Figure: Distribution of learnt weights after pruning small weights

Distribution of Learnt Weights





Figure: Distribution of learnt weights before pruning

Figure: Distribution of learnt weights after pruning small weights

- Of the 3 Billion weights, 97% are s.t, $W_{d',m'} \leq |0.01|,$ and hence non-discriminative
- · Storing them leads to large model sizes but no benefit in classification

DiSMEC - Distributed Sparse Machines for XMC

Require: Training data $\mathcal{T} = \{(\boldsymbol{x}_1, \boldsymbol{y}_1) \dots (\boldsymbol{x}_n, \boldsymbol{y}_n)\}$, input dimensionality d, label set $\{1 \dots m\}$, $B = \lfloor \frac{m}{1000} \rfloor + 1$ and Δ

Ensure: Learnt matrix $W_{d,m}$ in sparse format

- 1: Load single copy of input vectors $oldsymbol{X} = \{oldsymbol{x}_1 \dots oldsymbol{x}_n\}$ in the main memory
- 2: Load binary sign vectors $s_{m'} = \{+1, -1\}_{i=1}^n$ separately for each label
- 3: for $\{b = 0; b < B; b + +\}$ do 4: #pragma omp parallel for private(m') 5: for $\{m' = b \times 1000; m' \le (b + 1) \times 1000; m' + +\}$ do 6: Using $(\mathbf{X}, \mathbf{s}_{m'})$, train weight vector $\mathbf{w}_{m'}$ on a single core 7: Prune ambiguous weights in $\mathbf{w}_{m'}$ 8: return $\mathbf{W}_{d,1000}$ V Learnt matrix for a batch on one node
- 9: return $oldsymbol{W}_{d,m}$

> Learnt matrix from all the nodes

- Learns model for LSHTCWiki-325K in 6 hours on 400 cores
- Model size is 3GB due to pruning step

PDSparse - Primal Dual Sparsity for Extreme Classification ¹³

Primal Dual Sparsity in Max-margin Multi-label SVM Formulation

• For prediction score $z \in \mathbb{R}^m$ and true label $y \in \{1, \ldots, m\}$, $L(z, y) = \max_{m_- \in \{\mathcal{N}(y)\}, m_+ \in \{\mathcal{P}(y)\}} (1 + z_{m_-} - z_{m_+})_+$, where $\mathcal{N}(y)$ and $\mathcal{P}(y)$ are respectively the sets of negative and positive labels for that instance

¹³ Ian E.H. Yen, Xiangru Huang, Kai Zhong, Pradeep Ravikumar, and Inderjit S. Dhillon. PD-Sparse: A Primal and Dual Sparse Approach to Extreme Multiclass and Multilabel Classification. In *ICML*, 2016

PDSparse - Primal Dual Sparsity for Extreme Classification ¹³

Primal Dual Sparsity in Max-margin Multi-label SVM Formulation

- For prediction score $\boldsymbol{z} \in \mathbb{R}^m$ and true label $\boldsymbol{y} \in \{1, \ldots, m\}$, $L(\boldsymbol{z}, \boldsymbol{y}) = \max_{m_- \in \{\mathcal{N}(\boldsymbol{y})\}, m_+ \in \{\mathcal{P}(\boldsymbol{y})\}} (1 + z_{m_-} - z_{m_+})_+$, where $\mathcal{N}(\boldsymbol{y})$ and $\mathcal{P}(\boldsymbol{y})$ are respectively the sets of negative and positive labels for that instance
- $W^* \in \mathbb{R}^{d \times m}$ obtained by $\sum_{i=1}^n L(W^T x_i, y_i)$ is determined by the scores of *support labels* :

$$(m_{-}, m_{+}) = \arg\max_{m_{-} \in \{\mathcal{N}(\boldsymbol{y})\}, m_{+} \in \{\mathcal{P}(\boldsymbol{y})\}} (1 + z_{m_{-}} - z_{m_{+}})_{+}$$

that attain the above maximum

¹³ Ian E.H. Yen, Xiangru Huang, Kai Zhong, Pradeep Ravikumar, and Inderjit S. Dhillon. PD-Sparse: A Primal and Dual Sparse Approach to Extreme Multiclass and Multilabel Classification. In *ICML*, 2016

PDSparse - Primal Dual Sparsity for Extreme Classification ¹³

Primal Dual Sparsity in Max-margin Multi-label SVM Formulation

- For prediction score $\boldsymbol{z} \in \mathbb{R}^m$ and true label $\boldsymbol{y} \in \{1, \ldots, m\}$, $L(\boldsymbol{z}, \boldsymbol{y}) = \max_{m_- \in \{\mathcal{N}(\boldsymbol{y})\}, m_+ \in \{\mathcal{P}(\boldsymbol{y})\}} (1 + z_{m_-} - z_{m_+})_+$, where $\mathcal{N}(\boldsymbol{y})$ and $\mathcal{P}(\boldsymbol{y})$ are respectively the sets of negative and positive labels for that instance
- $W^* \in \mathbb{R}^{d \times m}$ obtained by $\sum_{i=1}^n L(W^T x_i, y_i)$ is determined by the scores of *support labels* :

$$(m_{-}, m_{+}) = \arg\max_{m_{-} \in \{\mathcal{N}(\boldsymbol{y})\}, m_{+} \in \{\mathcal{P}(\boldsymbol{y})\}} (1 + z_{m_{-}} - z_{m_{+}})_{+}$$

that attain the above maximum

• On the dual side, since the max is attained for very few label pairs (m_-, m_+) , the dual variable is non-zero very infrequently i.e. $nnz(\alpha_i) << m$

¹³ Ian E.H. Yen, Xiangru Huang, Kai Zhong, Pradeep Ravikumar, and Inderjit S. Dhillon. PD-Sparse: A Primal and Dual Sparse Approach to Extreme Multiclass and Multilabel Classification. In *ICML*, 2016

PDSparse - Primal Dual Sparsity for Extreme Classification

Primal Dual Sparsity in Max-margin Multi-label SVM Formulation

• On the primal side, this is achieved by ℓ_1 regularization

$$oldsymbol{W}^* \in \operatorname*{arg\,min}_{oldsymbol{W}} \lambda \sum_{m'=1}^m ||oldsymbol{w}_{m'}||_1 + \sum_{i=1}^n L(oldsymbol{W}^T oldsymbol{x}_i, oldsymbol{y}_i)$$

then it satisfies $d \times m_w = nnz(\mathbf{W}^*) \le nnz(\mathbf{A}^*) = n \times m_a$, where m_w are average number of active labels per feature in the optimal weight matrix \mathbf{W}^* .

PDSparse - Primal Dual Sparsity for Extreme Classification

Primal Dual Sparsity in Max-margin Multi-label SVM Formulation

• For ease of optimization, the following *strongly convex* elastic net regularized problem is solved

$$\underset{\boldsymbol{W}}{\operatorname{arg\,min}} \sum_{m'=1}^{m} ||\boldsymbol{w}_{m'}||_{2}^{2} + \lambda \sum_{m'=1}^{m} ||\boldsymbol{w}_{m'}||_{1} + \sum_{i=1}^{n} L(\boldsymbol{W}^{T}\boldsymbol{x}_{i}, \boldsymbol{y}_{i})$$

which is empirically observed to give a similar sparsity level as the ℓ_1 regularized problem

PDSparse - Primal Dual Sparsity for Extreme Classification

Simultaneous primal dual updates



Figure: Simultaneous Primal Dual Updates in PDSparse

Dataset and Evaluation Metrics for Comparison

Datasets taken from Extreme Classification repository

| Dataset | # Training | # Test | # Categories | # Features | | |
|----------------|------------|-----------|--------------|------------|--------|------|
| | | | | | APpL | ALpP |
| Amazon-13K | 1,186,239 | 306,782 | 13,330 | 203,882 | 448.5 | 5.04 |
| Amazon-14K | 4,398,050 | 1,099,725 | 14,588 | 597,540 | 1330.1 | 3.53 |
| Wikipedia-31K | 14,146 | 6,616 | 30,938 | 101,938 | 8.5 | 18.6 |
| Delicious-200K | 196,606 | 100,095 | 205,443 | 1,123,497 | 72.3 | 75.5 |
| WikiLSHTC-325K | 1,778,351 | 587,084 | 325,056 | 1,617,899 | 17.4 | 3.2 |
| Wikipedia-500K | 1,813,391 | 783,743 | 501,070 | 2,381,304 | 24.7 | 4.7 |
| Amazon-670K | 490,499 | 153,025 | 670,091 | 135909 | 3.9 | 5.4 |

Table: Multi-label datasets taken from the Extreme Classification Repository. APpL and ALpP represent average points per label and average labels per point.

For true label vector $oldsymbol{y}$, and predicted vector $\hat{oldsymbol{y}}$:

$$\mathsf{precision} @\mathsf{k} := \frac{1}{k} \sum_{l \in rank_k(\hat{\pmb{y}})} \pmb{y}_l \hspace{3mm} ; \hspace{3mm} \mathsf{nDCG}@\mathsf{k} := \frac{\mathsf{DCG}@\mathsf{k}}{\sum_{l=1}^{\min(k, ||\pmb{y}||_0)} \frac{1}{\log(l+1)}}$$

where DCG@k := $\frac{1}{k} \sum_{l \in rank_k(\hat{y})} \frac{y_l}{\log(l+1)}$, and $rank_k(y)$ returns the k largest indices of y ranked in descending order, and $||y||_0$ returns the 0-norm of the true-label vector.

Comparison among state-of-the-art methods on Precision@k metric

| Dataset | Proposed approach Embedding based approaches | | | ed approaches | Tree based approaches | | | Sparsity inducing approaches | |
|----------------|--|-----------|------|---------------|-----------------------|---------|------|------------------------------|--------|
| Dataset | DiSMEC | SLEEC | LEML | RobustXML | Fast-XML | LPSR-NB | PLT | PD-Sparse | L1-SVM |
| Amazon-13K | | | | | | | | | |
| P@1 | 93.4 | 90.4 | 78.2 | 88.4 | 92.9 | 75.1 | 91.4 | 91.1 | 91.8 |
| P@3 | 79.1 | 75.8 | 65.4 | 74.6 | 77.5 | 60.2 | 75.8 | 76.4 | 77.8 |
| P@5 | 64.1 | 61.3 | 55.7 | 60.6 | 62.5 | 57.3 | 61.0 | 63.0.8 | 62.9 |
| Amazon-14K | | | | | | | | | |
| P@1 | 91.0 | 80.3 | 75.2 | 83.2 | 90.3 | 74.2 | 86.4 | 88.4 | 88.2 |
| P@3 | 70.3 | 67.2 | 62.5 | 66.4 | 70.1 | 55.7 | 65.2 | 68.1 | 67.6 |
| P@5 | 55.9 | 50.6 | 40.8 | 52.3 | 55.4 | 44.3 | 50.7 | 50.5 | 51.2 |
| Wikipedia-31K | | | | | | | | | |
| P@1 | 85.2 | 85.5 | 73.5 | 85.5 | 82.5 | 72.7 | 84.3 | 73.8 | 83.2 |
| P@3 | 74.6 | 73.6 | 62.3 | 74.0 | 66.6 | 58.5 | 72.3 | 60.9 | 72.1 |
| P@5 | 65.9 | 63.1 | 54.3 | 63.8 | 56.7 | 49.4 | 62.7 | 50.4 | 63.7 |
| Delicious-200k | | | | | | | | | |
| P@1 | 45.5 | 47.0 40.3 | | 45.0 | 42.8 | 18.6 | 45.3 | 41.2 | 42.1 |
| P@3 | 38.7 | 41.6 | 37.7 | 40.0 | 38.7 | 15.4 | 38.9 | 35.3 | 34.8 |
| P@5 | 35.5 | 38.8 | 36.6 | 38.0 | 36.3 | 14.0 | 35.8 | 31.2 | 30.4 |
| WikiLSHTC-325K | | | | | | | | | |
| P@1 | 64.4 | 55.5 | 19.8 | 53.5 | 49.3 | 27.4 | 45.6 | 58.2 | 60.6 |
| P@3 | 42.5 | 33.8 | 11.4 | 31.8 | 32.7 | 16.4 | 29.1 | 36.3 | 38.6 |
| P@5 | 31.5 | 24.0 | 8.4 | 29.9 | 24.0 | 12.0 | 21.9 | 28.7 | 28.5 |
| Wiki-500K | | | | | | | | | |
| P@1 | 70.2 48.2 4 | | 41.3 | - | 54.1 | 38.2 | 51.5 | - | 65.3 |
| P@3 | 50.6 | 50.6 29.4 | | - | 35.5 | 29.3 | 35.7 | - | 46.1 |
| P@5 | 39.7 | 21.2 | 19.8 | - | 26.2 | 18.7 | 27.7 | - | 35.3 |
| Amazon-670K | | | | | | | | | |
| P@1 | 44.7 | 35.0 | 8.1 | 31.0 | 33.3 | 28.6 | 36.6 | - | 39.8 |
| P@3 | 39.7 | 31.2 | 6.8 | 28.0 | 29.3 | 24.9 | 32.1 | - | 34.3 |
| P@5 | 36.1 | 28.5 | 6.0 | 24.0 | 26.1 | 22.3 | 28.8 | - | 30.1 |

Table: Comparison of Precision@k for k=1,3 and 5

Comparison among methods on nDCG@k metric



Figure: nDCG for Amazon-670K



Figure: nDCG for Delicious-200K



Figure: nDCG for WikiLSHTC-325K



Figure: nDCG for Wiki-30K

Impact of Δ parameter



Figure: Variation of Precision@k for WikiLSHTC-325K with Δ

Figure: Variation of Model size for WikiLSHTC-325K with Δ

Label Embedding versus One-vs-Rest - What to choose when?

Algebraic graph theoretic view-point

- Let A(G) be adjacency matrix for the label co-occurrence graph
- Let ${\cal D}({\cal G})$ be degree matrix representing degrees of labels

The Laplacian of the graph G, is given by L(G) = D(G) - A(G).

Theorem

¹⁴ Given the training data \mathcal{T} , let L(G) be defined as above. Let $\lambda_1(G), \ldots, \lambda_L(G)$ be the eigen-values of L(G), then : (a) $\lambda_\ell(G) \ge 0 \quad \forall \ell$, and $\lambda_1(G) = 0$ (b) The multiplicity of 0 as an eigen value gives the number of connected components of G, $\implies \lambda_2(G) > 0$ if and only if G is connected. (c) $\lambda_2(G) \le \nu(G) \le \eta(G)$, where $\lambda_2(G), \nu(G)$ and $\eta(G)$ are respectively the algebraic, edge and vertex connectivities.

¹⁴ Miroslav Fiedler. Algebraic connectivity of graphs. Czechoslovak mathematical journal, 23(2):298–305, 1973

Label Embedding vs One-vs-Rest - What to choose When?

| Dataset | # Training | # Categories | | | Algebraic | # Connected | Better |
|-----------------|------------|--------------|--------|------|------------------------------|-------------|-------------|
| | | | APpL | ALpP | Connectivity, $\lambda_2(G)$ | Components | Performance |
| Mediamill | 30993 | 101 | 1902.1 | 4.4 | 0.46 | 1 | SLEEC |
| Bibtex | 4880 | 159 | 111.7 | 2.4 | 0.05 | 1 | SLEEC |
| Delicious-small | 12,920 | 983 | 311.6 | 19.3 | 0.3 | 1 | SLEEC |
| EUR-Lex | 15,539 | 3,993 | 25.7 | 5.3 | 0.22 | 1 | DiSMEC |
| Wikipedia-31K | 14,146 | 30,938 | 8.5 | 18.6 | 0.4 | 1 | Comparable |
| WikiLSHTC-325K | 1,778,351 | 325,056 | 17.4 | 3.2 | 0.002 | 740 | DiSMEC |
| Wiki-500K | 1,813,391 | 501,070 | 24.7 | 4.7 | 0.001 | 370 | DiSMEC |
| Amazon-670K | 490,499 | 670,091 | 3.9 | 5.4 | 0.0001 | 9,566 | DiSMEC |

Table: Multi-label datasets from XMC repository. APpL and ALpP represent average points per label and average labels per point respectively

• For smaller datasets, the algebraic connectivity is much higher and label embedding based methods such as SLEEC work relatively well in this regime

Label Embedding vs One-vs-Rest - What to choose When?

| Dataset | # Training | # Categories | | | Algebraic | # Connected | Better |
|-----------------|------------|--------------|--------|------|------------------------------|-------------|-------------|
| | | | APpL | ALpP | Connectivity, $\lambda_2(G)$ | Components | Performance |
| Mediamill | 30993 | 101 | 1902.1 | 4.4 | 0.46 | 1 | SLEEC |
| Bibtex | 4880 | 159 | 111.7 | 2.4 | 0.05 | 1 | SLEEC |
| Delicious-small | 12,920 | 983 | 311.6 | 19.3 | 0.3 | 1 | SLEEC |
| EUR-Lex | 15,539 | 3,993 | 25.7 | 5.3 | 0.22 | 1 | DiSMEC |
| Wikipedia-31K | 14,146 | 30,938 | 8.5 | 18.6 | 0.4 | 1 | Comparable |
| WikiLSHTC-325K | 1,778,351 | 325,056 | 17.4 | 3.2 | 0.002 | 740 | DiSMEC |
| Wiki-500K | 1,813,391 | 501,070 | 24.7 | 4.7 | 0.001 | 370 | DiSMEC |
| Amazon-670K | 490,499 | 670,091 | 3.9 | 5.4 | 0.0001 | 9,566 | DiSMEC |

Table: Multi-label datasets from XMC repository. APpL and ALpP represent average points per label and average labels per point respectively

- For smaller datasets, the algebraic connectivity is much higher and label embedding based methods such as SLEEC work relatively well in this regime
- For much larger datasets, the algebraic connectivity is close to 0 and one-vs-rest approach such as DiSMEC works better in this regime

Tree-based models

Tree-based models

- Decision trees
- Label trees

• Decision trees:

¹⁵ Anna Choromanska and John Langford. Logarithmic time online multiclass prediction. In NIPS 29, 2015

¹⁶Yashoteja Prabhu and Manik Varma. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In KDD, pages 263–272. ACM, 2014

- Decision trees:
 - Partition of the feature space to small subregions:



¹⁵ Anna Choromanska and John Langford. Logarithmic time online multiclass prediction. In NIPS 29, 2015

¹⁶Yashoteja Prabhu and Manik Varma. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In KDD, pages 263–272. ACM, 2014

- Decision trees:
 - Partition of the feature space to small subregions:



• Fast prediction: logarithmic in n

¹⁵ Anna Choromanska and John Langford. Logarithmic time online multiclass prediction. In NIPS 29, 2015

¹⁶Yashoteja Prabhu and Manik Varma. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In KDD, pages 263–272. ACM, 2014

- Decision trees:
 - Partition of the feature space to small subregions:



- ▶ Fast prediction: logarithmic in n
- Training can be expensive: computation of split criterion

¹⁵ Anna Choromanska and John Langford. Logarithmic time online multiclass prediction. In NIPS 29, 2015

¹⁶Yashoteja Prabhu and Manik Varma. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In KDD, pages 263–272. ACM, 2014

- Decision trees:
 - ► Partition of the feature space to small subregions:



- ▶ Fast prediction: logarithmic in n
- Training can be expensive: computation of split criterion
- ► Two new algorithms: LomTree¹⁵ and FastXML¹⁶

¹⁵ Anna Choromanska and John Langford. Logarithmic time online multiclass prediction. In NIPS 29, 2015

¹⁶Yashoteja Prabhu and Manik Varma. FastXML: A fast, accurate and stable tree-classifier for extreme multi-label learning. In KDD, pages 263–272. ACM, 2014
FastXML

- Uses an ensemble of standard decision trees
- Sparse linear classifiers trained in internal nodes
- Very efficient training procedure
- Empirical distributions in leaves
- A test example passes one path from the root to a leaf



FastXML

- Uses an ensemble of standard decision trees
- Sparse linear classifiers trained in internal nodes
- Very efficient training procedure
- Empirical distributions in leaves
- A test example passes one path from the root to a leaf



Optimization in FastXML

• In each internal node FastXML solves:

$$\begin{array}{ll} \min & \|\boldsymbol{w}\|_1 + \sum_{i=1}^n C_{\delta}(\delta_i) \log(1 + \exp(-\delta_i \boldsymbol{w}^\top \boldsymbol{x}) \\ & -C_r \sum_{i=1}^n \frac{1}{2} (1 + \delta_i) \mathrm{NDCG}@m(\boldsymbol{r}^+, \boldsymbol{y}_i) \\ & -C_r \sum_{i=1}^n \frac{1}{2} (1 - \delta_i) \mathrm{NDCG}@m(\boldsymbol{r}^-, \boldsymbol{y}_i) \\ & \text{w.r.t.} & \boldsymbol{w} \in \mathbb{R}^d, \boldsymbol{\delta} \in \{-1, 1\}^n, \boldsymbol{r}^+, \boldsymbol{r}^- \in \Pi(1, m) \end{array}$$

Optimization in FastXML

• In each internal node FastXML solves:

$$\begin{array}{ll} \min & \|\boldsymbol{w}\|_{1} + \sum_{i=1}^{n} C_{\delta}(\delta_{i}) \log(1 + \exp(-\delta_{i} \boldsymbol{w}^{\top} \boldsymbol{x}) \\ & -C_{r} \sum_{i=1}^{n} \frac{1}{2} (1 + \delta_{i}) \text{NDCG}@m(\boldsymbol{r}^{+}, \boldsymbol{y}_{i}) \\ & -C_{r} \sum_{i=1}^{n} \frac{1}{2} (1 - \delta_{i}) \text{NDCG}@m(\boldsymbol{r}^{-}, \boldsymbol{y}_{i}) \\ \text{w.r.t.} & \boldsymbol{w} \in \mathbb{R}^{d}, \boldsymbol{\delta} \in \{-1, 1\}^{n}, \boldsymbol{r}^{+}, \boldsymbol{r}^{-} \in \Pi(1, m) \\ \\ \\ \text{linear split} & \text{label ranking in positive} \\ & \text{and negative partition} \end{array}$$

Optimization in FastXML

• In each internal node FastXML solves:

$$\begin{array}{ll} \min & \|\boldsymbol{w}\|_{1} + \sum_{i=1}^{n} C_{\delta}(\delta_{i}) \log(1 + \exp(-\delta_{i}\boldsymbol{w}^{\top}\boldsymbol{x}) \\ & \text{1. Bernoulli sampling of } \delta \\ & \text{2. Optimization of } \boldsymbol{r}^{\pm} \\ & \text{3. Optimization of } \delta \\ & \text{4. Optimization of } \boldsymbol{w} \\ & \text{5. Repeat 2-4} \end{array} \qquad - C_{r} \sum_{i=1}^{n} \frac{1}{2}(1 + \delta_{i}) \text{NDCG}@m(\boldsymbol{r}^{+}, \boldsymbol{y}_{i}) \\ & - C_{r} \sum_{i=1}^{n} \frac{1}{2}(1 - \delta_{i}) \text{NDCG}@m(\boldsymbol{r}^{-}, \boldsymbol{y}_{i}) \\ & \text{w.r.t.} \qquad \boldsymbol{w} \in \mathbb{R}^{d}, \boldsymbol{\delta} \in \{-1, 1\}^{n}, \boldsymbol{r}^{+}, \boldsymbol{r}_{-} \in \Pi(1, m) \\ & \text{linear split} \end{array}$$







¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA



Sampling of δ

¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA





¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA



Optimization of δ



¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA





¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA



Optimization of δ



¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA





¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA



Optimization of \boldsymbol{w}



¹⁷https://www.youtube.com/watch?v=1X71fTx1LKA

- Label trees:
 - ► Organize classifiers in a tree structure (one leaf ⇔ one label):



¹⁸ S. Bengio, J. Weston, and D. Grangier. Label embedding trees for large multi-class tasks. In NIPS, pages 163–171. Curran Associates, Inc., 2010

¹⁹ J. Fox. Applied regression analysis, linear models, and related methods. Sage, 1997

E. Frank and S. Kramer. Ensembles of nested dichotomies for multi-class problems. In ICML, 2004

²⁰ A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In UAI, pages 51–58, 2009

²¹ Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In AISTATS, pages 246–252, 2005

²² Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. CoRR, abs/1607.01759, 2016

²³ K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In ICML, pages 279–286. Omnipress, 2010 756 (10)

- Label trees:
 - ► Organize classifiers in a tree structure (one leaf ⇔ one label):



▶ Fast prediction: almost logarithmic in m

¹⁸ S. Bengio, J. Weston, and D. Grangier. Label embedding trees for large multi-class tasks. In NIPS, pages 163–171. Curran Associates, Inc., 2010

¹⁹ J. Fox. Applied regression analysis, linear models, and related methods. Sage, 1997 E. Frank and S. Kramer. Ensembles of nested dichotomies for multi-class problems. In ICML, 2004

²⁰ A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In UAI, pages 51–58, 2009

²¹ Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In AISTATS, pages 246–252, 2005

²² Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. CoRR, abs/1607.01759, 2016

²³ K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In ICML, pages 279–286. Omnipress, 2010 756 (10)

- Label trees:
 - ► Organize classifiers in a tree structure (one leaf ⇔ one label):



- Fast prediction: almost logarithmic in m
- Algorithms: Label embedding trees,¹⁸ Nested dichotomies,¹⁹ Conditional probability trees,²⁰ Hierarchical softmax,²¹ FastText,²² Probabilistic classifier chains²³

¹⁸ S. Bengio, J. Weston, and D. Grangier. Label embedding trees for large multi-class tasks. In NIPS, pages 163–171. Curran Associates, Inc., 2010

¹⁹ J. Fox. Applied regression analysis, linear models, and related methods. Sage, 1997

E. Frank and S. Kramer. Ensembles of nested dichotomies for multi-class problems. In ICML, 2004

²⁰ A. Beygelzimer, J. Langford, Y. Lifshits, G. B. Sorkin, and A. L. Strehl. Conditional probability tree estimation analysis and algorithms. In UAI, pages 51–58, 2009

²¹ Frederic Morin and Yoshua Bengio. Hierarchical probabilistic neural network language model. In AISTATS, pages 246–252, 2005

²² Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. CoRR, abs/1607.01759, 2016

²³ K. Dembczyński, W. Cheng, and E. Hüllermeier. Bayes optimal multilabel classification via probabilistic classifier chains. In ICML, pages 279–286. Omnipress, 2010

Probabilistic classifier trees = Hierarchical softmax

Nested dichotomies, Hierarchical softmax, Conditional probability trees, Probabilistic classifier chains

↓ Probability classifier trees²⁴ ↓ Hierarchical softmax

²⁴ Krzysztof Dembczyński, Wojciech Kotłowski, Willem Waegeman, Róbert Busa-Fekete, and Eyke Hüllermeier. Consistency of probabilistic classifier trees. In ECMLPKDD. Springer, 2016



• Encode the labels by a **prefix code** (⇒ tree structure)



• Each label y coded by $\boldsymbol{z} = (z_1, \dots, z_l) \in \mathcal{C}$



- Each label y coded by $\boldsymbol{z} = (z_1, \dots, z_l) \in \mathcal{C}$
- An internal node identified by a partial code $oldsymbol{z}^i = (z_1,\ldots,z_i)$



- Each label y coded by $oldsymbol{z} = (z_1, \dots, z_l) \in \mathcal{C}$
- An internal node identified by a **partial** code $\boldsymbol{z}^i = (z_1, \dots, z_i)$
- The code does not have to be binary ⇒ b-ary trees



- Each label y coded by $oldsymbol{z} = (z_1, \dots, z_l) \in \mathcal{C}$
- An internal node identified by a partial code $oldsymbol{z}^i = (z_1,\ldots,z_i)$
- The code does **not** have to be binary \Rightarrow *b*-ary trees
- Different **structures** possible: random tree, Huffman tree, trained structure

• HSM estimates $\mathbf{P}(y \,|\, \boldsymbol{x})$ by following a **path** from the root to a leaf:



- Training: separate learning problems in the internal nodes
- Prediction: depth first search/beam search

• HSM estimates $\mathbf{P}(y | \boldsymbol{x})$ by following a **path** from the root to a leaf:



- Training: separate learning problems in the internal nodes
- Prediction: depth first search/beam search

HSM for MLC

• Pick-one-label heuristic used, for example, in FastText:

$$\eta'_j(\boldsymbol{x}) = \mathbf{P}'(y_j = 1 \,|\, \boldsymbol{x}) = \sum_{\boldsymbol{y} \in \mathcal{Y}} y_j \frac{\mathbf{P}(\boldsymbol{y} \,|\, \boldsymbol{x})}{\sum_{j'=1}^m y_{j'}}$$

HSM for MLC

• Pick-one-label heuristic used, for example, in FastText:

$$\eta_j'(\boldsymbol{x}) = \mathbf{P}'(y_j = 1 \,|\, \boldsymbol{x}) = \sum_{\boldsymbol{y} \in \mathcal{Y}} y_j \frac{\mathbf{P}(\boldsymbol{y} \,|\, \boldsymbol{x})}{\sum_{j'=1}^m y_{j'}}$$

• Theorem: inconsistent for label-wise logistic loss and precision@k

| labels $oldsymbol{y}$ | probability $\mathbf{P}(oldsymbol{y} oldsymbol{x})$ | True marg. probs | P-o-I marg. probs |
|-----------------------|---|---|---|
| | 0.15 0.10 0.25 0.30 0.20 | $egin{aligned} &\eta_1(m{x})=0.4\ &\eta_2(m{x})=0.35\ &\eta_3(m{x})=0.3\ &\eta_4(m{x})=0.2 \end{aligned}$ | $egin{aligned} &\eta_3'(m{x})=0.3 \ &\eta_1'(m{x})=0.275 \ &\eta_2'(m{x})=0.225 \ &\eta_4'(m{x})=0.2 \end{aligned}$ |

HSM for MLC

• Pick-one-label heuristic used, for example, in FastText:

$$\eta_j'(\boldsymbol{x}) = \mathbf{P}'(y_j = 1 \,|\, \boldsymbol{x}) = \sum_{\boldsymbol{y} \in \mathcal{Y}} y_j \frac{\mathbf{P}(\boldsymbol{y} \,|\, \boldsymbol{x})}{\sum_{j'=1}^m y_{j'}}$$

• Theorem: inconsistent for label-wise logistic loss and precision@k

| labels $m{y}$ | probability $\mathbf{P}(oldsymbol{y} oldsymbol{x})$ | True marg. probs | P-o-I marg. probs |
|---------------|---|--|---|
| | 0.15 0.10 0.25 0.30 0.20 | $egin{aligned} &\eta_1(m{x})=0.4 \ &\eta_2(m{x})=0.35 \ &\eta_3(m{x})=0.3 \ &\eta_4(m{x})=0.2 \end{aligned}$ | $egin{aligned} &\eta_3'(m{x})=0.3 \ &\eta_1'(m{x})=0.275 \ &\eta_2'(m{x})=0.225 \ &\eta_4'(m{x})=0.2 \end{aligned}$ |

• **Theorem**: consistent for precision@k for independent labels

• Similar tree structure and encoding of $y_j = 1$ by $\boldsymbol{z} = (1, z_1, \dots, z_l)$



²⁵ K. Jasinska, K. Dembczynski, R. Busa-Fekete, K. Pfannschmidt, T. Klerx, and E. Hüllermeier. Extreme F-measure maximization using sparse probability estimates. In *ICML*, 2016

• Similar tree structure and encoding of $y_j = 1$ by $\boldsymbol{z} = (1, z_1, \dots, z_l)$



• Marginal probabilities $\eta_j(x)$ obtained by:

$$\eta_j(\boldsymbol{x}) = \mathbf{P}(\boldsymbol{z} \mid \boldsymbol{x}) = \prod_{i=0}^{n} \mathbf{P}(z_i \mid \boldsymbol{z}^{i-1}, \boldsymbol{x})$$

²⁵ K. Jasinska, K. Dembczynski, R. Busa-Fekete, K. Pfannschmidt, T. Klerx, and E. Hüllermeier. Extreme F-measure maximization using sparse probability estimates. In *ICML*, 2016

• Similar tree structure and encoding of $y_j = 1$ by $\boldsymbol{z} = (1, z_1, \dots, z_l)$



• Marginal probabilities $\eta_j(x)$ obtained by:

$$\eta_j(\boldsymbol{x}) = \mathbf{P}(\boldsymbol{z} \mid \boldsymbol{x}) = \prod_{i=0}^{r} \mathbf{P}(z_i \mid \boldsymbol{z}^{i-1}, \boldsymbol{x})$$

• $z^i \Leftrightarrow$ at least one positive label in the corresponding subtree

²⁵ K. Jasinska, K. Dembczynski, R. Busa-Fekete, K. Pfannschmidt, T. Klerx, and E. Hüllermeier. Extreme F-measure maximization using sparse probability estimates. In *ICML*, 2016

• Similar tree structure and encoding of $y_j = 1$ by $\boldsymbol{z} = (1, z_1, \dots, z_l)$



• Marginal probabilities $\eta_j(x)$ obtained by:

$$\eta_j(\boldsymbol{x}) = \mathbf{P}(\boldsymbol{z} \mid \boldsymbol{x}) = \prod_{i=0}^{r} \mathbf{P}(z_i \mid \boldsymbol{z}^{i-1}, \boldsymbol{x})$$

• $\boldsymbol{z}^i \Leftrightarrow \operatorname{\mathsf{at}}$ least one positive label in the corresponding subtree

• $\sum_{z_t} \mathbf{P}(z_t \,|\, \boldsymbol{z}^{t-1}) \geq 1 \Rightarrow$ separate classifiers in all nodes of the tree

²⁵ K. Jasinska, K. Dembczynski, R. Busa-Fekete, K. Pfannschmidt, T. Klerx, and E. Hüllermeier. Extreme F-measure maximization using sparse probability estimates. In *ICML*, 2016

• Training:

- independent training of all node classifiers
- reduced complexity by the conditions used in the nodes
- batch or online learning of node classifiers
- sparse representation: small number of active features in lower nodes, feature hashing
- dense representation: hidden representation of features (strong compression)

• Training:

- independent training of all node classifiers
- reduced complexity by the conditions used in the nodes
- batch or online learning of node classifiers
- sparse representation: small number of active features in lower nodes, feature hashing
- dense representation: hidden representation of features (strong compression)

• Theoretical guarantees

▶ Regret bounds for the absolute error of marginal probability estimates and precision@k

• Training:

- independent training of all node classifiers
- reduced complexity by the conditions used in the nodes
- batch or online learning of node classifiers
- sparse representation: small number of active features in lower nodes, feature hashing
- dense representation: hidden representation of features (strong compression)

• Theoretical guarantees

▶ Regret bounds for the absolute error of marginal probability estimates and precision@k

• Tree structure:

► Random, complete-sorted tree, Huffman tree, trained (e.g., hierarchical clustering of labels), online training

• Training:

- independent training of all node classifiers
- reduced complexity by the conditions used in the nodes
- batch or online learning of node classifiers
- sparse representation: small number of active features in lower nodes, feature hashing
- dense representation: hidden representation of features (strong compression)

• Theoretical guarantees

▶ Regret bounds for the absolute error of marginal probability estimates and precision@k

• Tree structure:

► Random, complete-sorted tree, Huffman tree, trained (e.g., hierarchical clustering of labels), online training

• Prediction:

depth first search/beam search
²⁶ Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In WWW. ACM, 2018

• HSM with pick-one-label

- online training
- b-ary trees
- ► Vowpal Wabbit: complete-sorted tree and feature hashing
- ► FastText: Huffman tree and hidden representation of features

²⁶Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In WWW. ACM, 2018

• HSM with pick-one-label

- online training
- b-ary trees
- ► Vowpal Wabbit: complete-sorted tree and feature hashing
- ► FastText: Huffman tree and hidden representation of features

• PLT

- online training
- ► b-ary trees
- ► Vowpal Wabbit: complete-sorted tree and feature hashing
- ► FastText: Huffman tree and hidden representation of features

²⁶Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In WWW. ACM, 2018

• HSM with pick-one-label

- online training
- b-ary trees
- Vowpal Wabbit: complete-sorted tree and feature hashing
- ► FastText: Huffman tree and hidden representation of features

• PLT

- online training
- ► b-ary trees
- ► Vowpal Wabbit: complete-sorted tree and feature hashing
- ► FastText: Huffman tree and hidden representation of features
- Parabel²⁶
 - PLT-like algorithm
 - Batch learning with squared hinge loss
 - Binary tree with higher arity of the last tree level
 - ► Tree structure obtained via 2-means++ algorithm

²⁶Y. Prabhu, A. Kag, S. Harsola, R. Agrawal, and M. Varma. Parabel: Partitioned label trees for extreme classification with application to dynamic search advertising. In WWW. ACM, 2018

- Results on the WikiLSHTC dataset
 - ▶ Stats: m = 325056, d = 1617899, $n_{\text{test}} = 587084$, $n_{\text{train}} = 1778351$
 - ► Results for precision@k:

| Method | P@1 |
|--------------------------|-----------------------|
| HSM-vw PLT-vw | 36.90 41.63 |
| FastText PLT-FastText | 41.28 41.58 |
| Parabel | 59.23 |

Decision trees vs. Label trees

| | Decision trees | Label trees |
|---------------------------------|------------------------|--------------------|
| tree structure | \checkmark | \checkmark |
| structure learning | \checkmark | \checkmark |
| batch learning | \checkmark | \checkmark |
| online learning | \checkmark | \checkmark |
| number of trees | ≥ 1 | ≥ 1 |
| number of leaves | O(n) | m |
| internal node models | cuts/linear | linear/any model |
| prediction | empirical distribution | score along a path |
| visited paths during prediction | 1 | several |
| sparse probability estimation | \checkmark | \checkmark |

Label filtering

• Classification of a test example in case of linear models can be formulated as:

$$j^* = \operatorname*{arg\,max}_{j\in\{1,...,m\}} \boldsymbol{w}_j^\top \boldsymbol{x},$$

²⁷ Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In PODS '01, pages 102–113. ACM, New York, NY, USA, 2001

• Classification of a test example in case of linear models can be formulated as:

$$j^* = \operatorname*{arg\,max}_{j\in\{1,\ldots,m\}} \boldsymbol{w}_j^\top \boldsymbol{x},$$

i.e., as the problem of maximum inner product search (MIPS)

• Naive solution is O(m)

²⁷ Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In PODS '01, pages 102–113. ACM, New York, NY, USA, 2001

• Classification of a test example in case of linear models can be formulated as:

$$j^* = \operatorname*{arg\,max}_{j\in\{1,\dots,m\}} \boldsymbol{w}_j^\top \boldsymbol{x},$$

- Naive solution is O(m)
- Optimal exact algorithm: the threshold algorithm ²⁷

²⁷ Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In PODS '01, pages 102–113. ACM, New York, NY, USA, 2001

• Classification of a test example in case of linear models can be formulated as:

$$j^* = \operatorname*{arg\,max}_{j\in\{1,\ldots,m\}} \boldsymbol{w}_j^\top \boldsymbol{x},$$

- Naive solution is O(m)
- Optimal exact algorithm: the threshold algorithm ²⁷
 - ► Does not perform well with a large number of features and labels

²⁷ Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In PODS '01, pages 102–113. ACM, New York, NY, USA, 2001

• Classification of a test example in case of linear models can be formulated as:

$$j^* = \operatorname*{arg\,max}_{j\in\{1,\dots,m\}} \boldsymbol{w}_j^\top \boldsymbol{x},$$

- Naive solution is O(m)
- Optimal exact algorithm: the threshold algorithm ²⁷
 - ► Does not perform well with a large number of features and labels
- There is a need for approximate solution

²⁷ Ronald Fagin, Amnon Lotem, and Moni Naor. Optimal aggregation algorithms for middleware. In PODS '01, pages 102–113. ACM, New York, NY, USA, 2001

MIPS vs. nearest neighbors

MIPS is similar to the nearest neighbor search under the square or cosine distance:

$$egin{array}{rcl} j^{*} &=& rgmin_{j\in\{1,...,m\}} \|m{w}_{j}-m{x}\|_{2}^{2} = rgmax_{j\in\{1,...,m\}} m{w}_{j}^{ op}m{x} - rac{\|m{w}_{j}\|_{2}^{2}}{2} \ j^{*} &=& rgmax_{j\in\{1,...,m\}} rac{m{w}_{j}^{ op}m{x}}{\|m{w}_{j}\|\|m{x}\|} = rgmax_{j\in\{1,...,m\}} rac{m{w}_{j}^{ op}m{x}}{\|m{w}_{j}\|} \end{array}$$

²⁸ A. Shrivastava and P. Li. Improved asymmetric locality sensitive hashing (ALSH) for maximum inner product search (mips). In UAI, 2015

MIPS vs. nearest neighbors

 MIPS is similar to the nearest neighbor search under the square or cosine distance:

$$j^{*} = \underset{j \in \{1,...,m\}}{\arg\min} \|\boldsymbol{w}_{j} - \boldsymbol{x}\|_{2}^{2} = \underset{j \in \{1,...,m\}}{\arg\max} \boldsymbol{w}_{j}^{\top} \boldsymbol{x} - \frac{\|\boldsymbol{w}_{j}\|_{2}^{2}}{2}$$
$$j^{*} = \underset{j \in \{1,...,m\}}{\arg\max} \frac{\boldsymbol{w}_{j}^{\top} \boldsymbol{x}}{\|\boldsymbol{w}_{j}\| \|\boldsymbol{x}\|} = \underset{j \in \{1,...,m\}}{\arg\max} \frac{\boldsymbol{w}_{j}^{\top} \boldsymbol{x}}{\|\boldsymbol{w}_{j}\|}$$

- Some tricks needs to be used to treat MIPS as nearest neighbor search $^{\rm 28}$

....

²⁸ A. Shrivastava and P. Li. Improved asymmetric locality sensitive hashing (ALSH) for maximum inner product search (mips). In UAI, 2015

• A well-studied problem with many existing algorithms:

²⁹ J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. ACM Transactions on Mathematical Software, 3(3):209–226, 1977

³⁰ Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In ACM Symposium on Theory of Computing, STOC '98, pages 604–613, New York, NY, USA, 1998. ACM

³¹ Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. CoRR, abs/1702.08734, 2017

³² Alex Auvolat, Hugo Larochelle, Sarath Chandar, Pascal Vincent, and Yoshua Bengio. Clustering is efficient for approximate maximum inner product search. CoRR, abs/1507.05910, 2015

- A well-studied problem with many existing algorithms:
 - ► For low-dimensional problems efficient tree-based structures exist²⁹

²⁹ J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. ACM Transactions on Mathematical Software, 3(3):209–226, 1977

³⁰ Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In ACM Symposium on Theory of Computing, STOC '98, pages 604–613, New York, NY, USA, 1998. ACM

³¹ Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. CoRR, abs/1702.08734, 2017

³² Alex Auvolat, Hugo Larochelle, Sarath Chandar, Pascal Vincent, and Yoshua Bengio. Clustering is efficient for approximate maximum inner product search. CoRR, abs/1507.05910, 2015

- A well-studied problem with many existing algorithms:
 - ► For low-dimensional problems efficient tree-based structures exist²⁹
 - Approximate nearest neighbor search via locality-sensitive hashing³⁰

²⁹ J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. ACM Transactions on Mathematical Software, 3(3):209–226, 1977

³⁰ Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In ACM Symposium on Theory of Computing, STOC '98, pages 604–613, New York, NY, USA, 1998. ACM

³¹ Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. CoRR, abs/1702.08734, 2017

³² Alex Auvolat, Hugo Larochelle, Sarath Chandar, Pascal Vincent, and Yoshua Bengio. Clustering is efficient for approximate maximum inner product search. CoRR, abs/1507.05910, 2015

- A well-studied problem with many existing algorithms:
 - ► For low-dimensional problems efficient tree-based structures exist²⁹
 - Approximate nearest neighbor search via locality-sensitive hashing³⁰
 - ► Fast search via k-means³¹ or hierarchical k-means³²

²⁹ J. H. Friedman, J. L. Bentley, and R. A. Finkel. An algorithm for finding best matches in logarithmic expected time. ACM Transactions on Mathematical Software, 3(3):209–226, 1977

³⁰ Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: Towards removing the curse of dimensionality. In ACM Symposium on Theory of Computing, STOC '98, pages 604–613, New York, NY, USA, 1998. ACM

³¹ Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. CoRR, abs/1702.08734, 2017

³² Alex Auvolat, Hugo Larochelle, Sarath Chandar, Pascal Vincent, and Yoshua Bengio. Clustering is efficient for approximate maximum inner product search. CoRR, abs/1507.05910, 2015

Agenda

- 1 Extreme classification: applications and challenges
- Algorithms
- **3 Live demonstration**

Live demonstration

- DataSets
 - ► Around a dozen datasets from Wikipedia, Amazon, Delicious
 - Datasets in raw and tf-idf format

³³ http://manikvarma.org/downloads/XC/XMLRepository.html

- DataSets
 - ► Around a dozen datasets from Wikipedia, Amazon, Delicious
 - Datasets in raw and tf-idf format
- Code
 - Around a dozen open source codes
 - State-of-the-art embedding, tree and linear methods

³³ http://manikvarma.org/downloads/XC/XMLRepository.html

- DataSets
 - ► Around a dozen datasets from Wikipedia, Amazon, Delicious
 - Datasets in raw and tf-idf format
- Code
 - Around a dozen open source codes
 - ► State-of-the-art embedding, tree and linear methods
- Comparisons
 - Comparison among state-of-the-art methods on various datasets
 - ► On Precision@k and nDCG@k and their propensity scored variants

³³ http://manikvarma.org/downloads/XC/XMLRepository.html

- DataSets
 - ► Around a dozen datasets from Wikipedia, Amazon, Delicious
 - Datasets in raw and tf-idf format
- Code
 - Around a dozen open source codes
 - ► State-of-the-art embedding, tree and linear methods
- Comparisons
 - Comparison among state-of-the-art methods on various datasets
 - ► On Precision@k and nDCG@k and their propensity scored variants
- Papers
 - ► Links to most of the recent papers in extreme classification

³³ http://manikvarma.org/downloads/XC/XMLRepository.html

• Extreme classification: #examples, #features, #labels

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,
 - smart 1-vs-All approaches,

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,
 - smart 1-vs-All approaches,
 - tree-based methods,

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,
 - smart 1-vs-All approaches,
 - tree-based methods,
 - ► label filtering/maximum inner product search.

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,
 - smart 1-vs-All approaches,
 - tree-based methods,
 - ► label filtering/maximum inner product search.

http://www.cs.put.poznan.pl/kdembczynski/ xmlc-tutorial-ecir-2018/

- Extreme classification: #examples, #features, #labels
- Complexity: time vs. space, training vs. validation vs. prediction
- Computational and stastistical challenges
- Different learning paradigms:
 - label embeddings,
 - smart 1-vs-All approaches,
 - tree-based methods,
 - ► label filtering/maximum inner product search.

http://www.cs.put.poznan.pl/kdembczynski/ xmlc-tutorial-ecir-2018/

Acknowledgments: Kalina Jasinska, Marek Wydmuch, Robert Busa-Fekete, Eyke Hüllermeier, Bernhard Schölkopf