



Extreme Classification

Tighter Bounds, Distributed Training, and new Algorithms

Marius Kloft (HU Berlin)

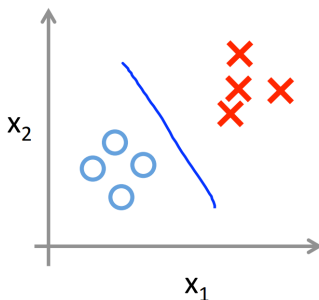
Krakow, Feb 16, 2017

- 1 Introduction
- 2 Distributed Algorithms
- 3 Theory
- 4 Learning Algorithms
- 5 Conclusion

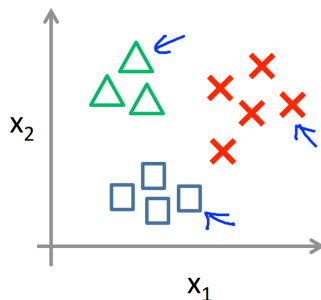
What is Multi-class Classification?

Multiclass classification is, given a data point x , decide on the class with which the data point is annotated.

Binary classification:



Multi-class classification:

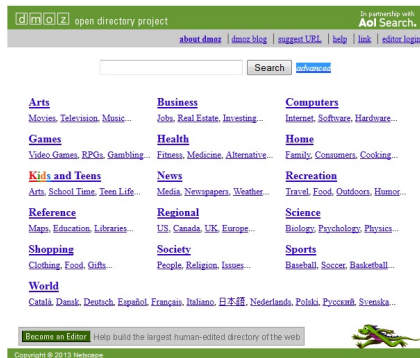


What is Extreme Classification?

Extreme classification is multi-class classification using an extremely large amount of classes.

Example 1

We are continuously monitoring the internet for new webpages, which we would like to categorize.



Example 2

We have data from an online biomedical bibliographic database that we want to index for quick access to clinicians.

NCBI Resources How To

PubMed.gov
U.S. National Library of Medicine
National Institutes of Health

Search: PubMed
blood

RSS Save search Adv

Display Settings: Summary, 20 per page, Sorted by Recently Added
Send to: Filter your results

Results: 1 to 20 of 2982326

1. [Toxic effects of Litsea elliptica Blume essential oil on Sprague-Dawley rats.](#)
Taib IS, Budin SB, Siti Nor Ain SM, Mohamed J, Lot S, Rajab NF, Hidayatulfathi O.
J Zhejiang Univ Sci B. 2009 Nov;10(11):813-9.
PMID: 19882755 [PubMed - in process]
[Related articles](#)

2. [Gasless laparoscopy for benign gynecological disease: wall-lifting system.](#)
Wang Y, Cui H, Zhao Y, Wang ZQ.
J Zhejiang Univ Sci B. 2009 Nov;10(11):805-12.

Choose Destination

☒ File ☐ Clipboard
☐ Collections ☐ E-mail
☐ Order

Download 2982326 items.

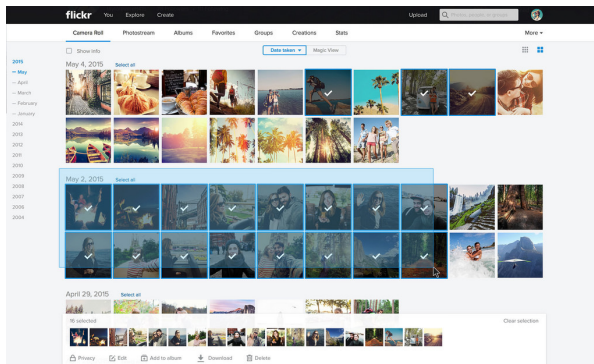
Format
MEDLINE

Sort by
Recently Added

Create File

Example 3

We are collecting data from an online feed of photographs that we would like to classify into image categories.



Example 4

We add new articles to an online encyclopedia and intend to predict the categories of the articles.

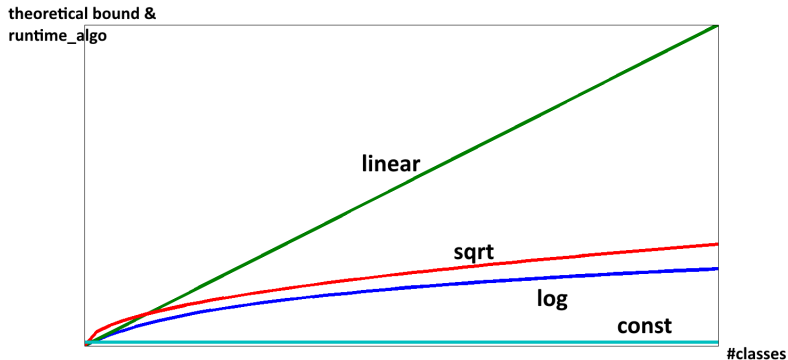


Need

Need for theory and algorithms for **extreme classification**.

How do **algorithms** and **bounds** scale
in **#classes**?

How do **algorithms** and **bounds** scale in **#classes**?

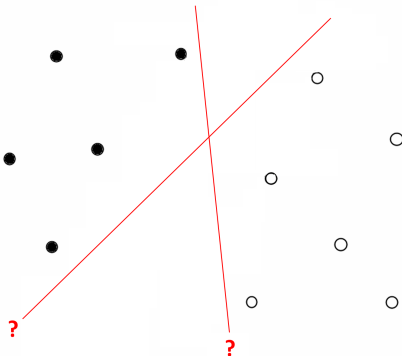


- 1 Introduction
- 2 Distributed Algorithms**
- 3 Theory
- 4 Learning Algorithms
- 5 Conclusion

Support Vector Machine (SVM) is a Popular Method for Binary Classification (Cortes and Vapnik, '95)

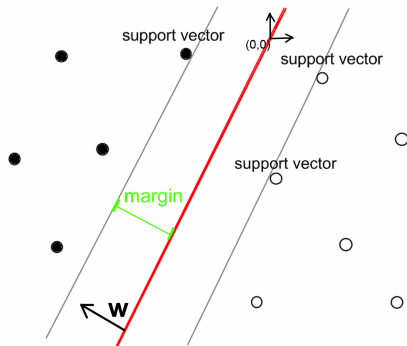
Core idea:

- Which hyperplane to take?



Support Vector Machine (SVM) is a Popular Method for Binary Classification

- ▶ Which hyperplane to take?
- ▶ **The one that separates the data with the largest margin**

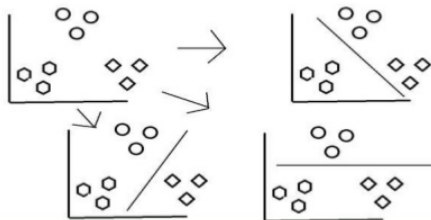


Popular Generalization to Multiple Classes: One-vs.-Rest SVM

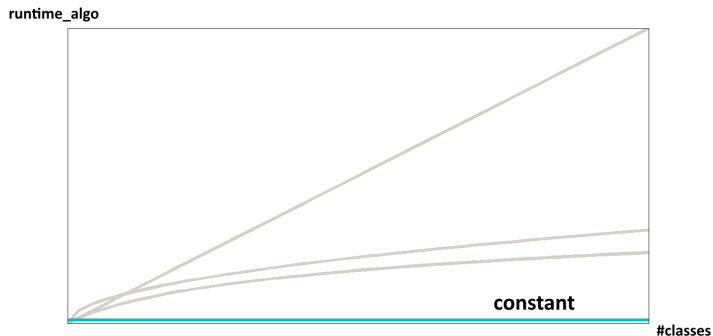
Put $\mathcal{C} := \#\text{classes}$.

One-vs.-rest SVM

- 1 For $c = 1..\mathcal{C}$
- 2 $\text{class1} := c, \text{class2} := \text{union}(\text{allOtherClasses})$
- 3 $w_c := \text{solutionOfSVM}(\text{class1}, \text{class2})$
- 4 end
- 5 Given a test point x , predict $c_{\text{predicted}} := \arg \max_c w_c^\top x$



Runtime of One-vs.-Rest



... assuming sufficient computational resources (`#classes` many computers)

Problem With One-vs.-Rest

:) training **can be parallelized** in the number of classes
(extreme classification!)

:(Is just a hack. One-vs.-Rest SVM is not built for multiple
classes (coupling of classes not exploited)!

There are “True” Multi-class SVMs, So-called **All-in-one** Multi-class SVMs

binary:

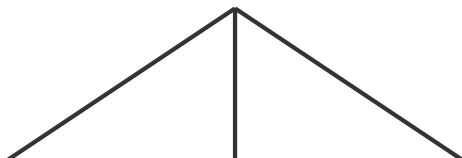
SVM

MC:

Lin, Lee, and
Wahba ('04)

Watkins and
Weston ('99)

Crammer and
Singer ('02)



There are “True” Multi-class SVMs, So-called **All-in-one** Multi-class SVMs

binary:

SVM

MC:

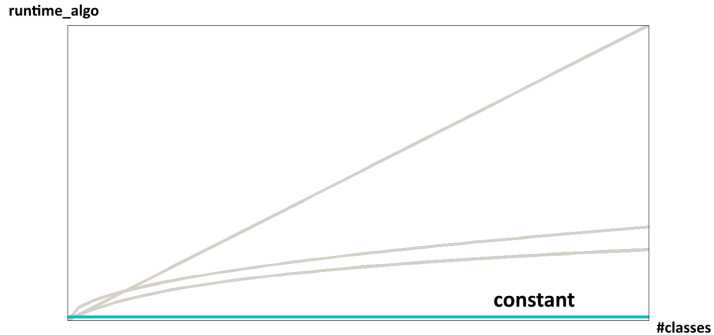
Lin, Lee, and
Wahba ('04)

Watkins and
Weston ('99)

Crammer and
Singer ('02)

Problem: State of the art solvers require a training time complexity of $\mathcal{O}(dn \cdot C)$, where $d = \text{dim}$, $n = \text{\#examples}$, and $C := \text{\#classes}$.

Aim: Develop algorithms where $\mathcal{O}(C)$ machines in **parallel** and in $\mathcal{O}(dn)$ runtime train all-in-one MC-SVMs.



⇒ same time complexity as one-vs.-rest,
yet more sophisticated algorithm

All-in-one SVMs

All of them have in common that they minimize a trade-off of a regularizer and a loss term:

$$\min_{w=(w_1, \dots, C)} \frac{1}{2} \sum_c \|w_c\|^2 + C * L(w, \text{data})$$

All-in-one SVMs

All of them have in common that they minimize a trade-off of a regularizer and a loss term:

$$\min_{w=(w_1, \dots, C)} \frac{1}{2} \sum_c \|w_c\|^2 + C * L(w, \text{data})$$

All Three MC-SVMs have:

$$\min_{w=(w_1,\dots,w_C)} \frac{1}{2} \sum_c \|w_c\|^2 + C* \dots$$

All Three MC-SVMs have:

$$\min_{w=(w_1,\dots,w_C)} \frac{1}{2} \sum_c \|w_c\|^2 + C* \dots$$

But they differ in the loss:

note: $l(x) := \max(0, 1 - x)$

CS:
$$\dots \sum_{i=1}^n \left[\max_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

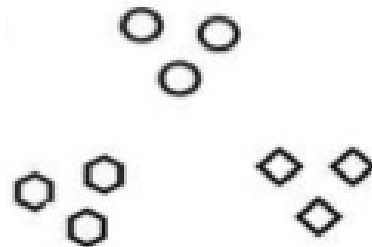
All Three MC-SVMs have:

$$\min_{w=(w_1, \dots, w_C)} \frac{1}{2} \sum_c \|w_c\|^2 + C * \dots$$

But they differ in the loss:

note: $l(x) := \max(0, 1 - x)$

CS:
$$\dots \sum_{i=1}^n \left[\max_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$



All Three MC-SVMs have:

$$\min_{w=(w_1,\dots,w_C)} \frac{1}{2} \sum_c \|w_c\|^2 + C * \dots$$

But they differ in the loss:

note: $l(x) := \max(0, 1 - x)$

$$\text{CS:} \quad \dots \sum_{i=1}^n \left[\max_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

$$\text{WW:} \quad \dots \sum_{i=1}^n \left[\sum_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

All Three MC-SVMs have:

$$\min_{w=(w_1,\dots,w_C)} \frac{1}{2} \sum_c \|w_c\|^2 + C * \dots$$

But they differ in the loss:

note: $l(x) := \max(0, 1 - x)$

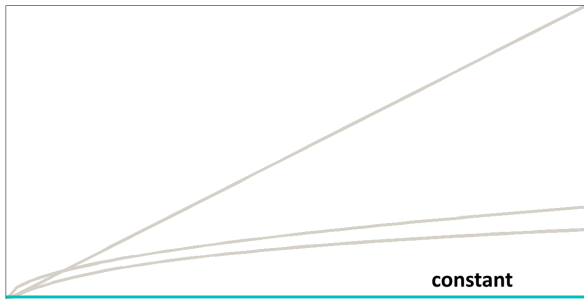
$$\text{CS:} \quad \dots \sum_{i=1}^n \left[\max_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

$$\text{WW:} \quad \dots \sum_{i=1}^n \left[\sum_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

$$\text{LLW:} \quad \dots \sum_{i=1}^n \left[\sum_{c \neq y_i} l(-w_c^T x_i) \right], \text{ s.t. } \sum_c w_c = 0$$

Can we solve these all-in-one MC-SVMs in parallel?

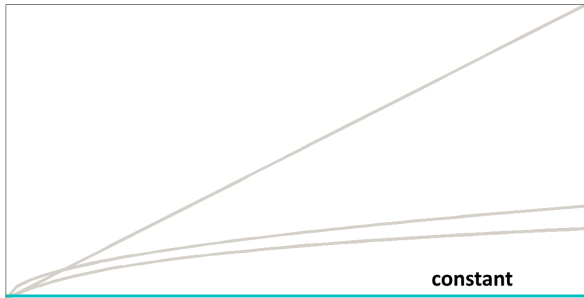
runtime_algo



#classes

Can we solve these all-in-one MC-SVMs in parallel?

runtime_algo



#classes

Let's look at **Lee, Lin, and Wahba (LLW)** first.

This is the LLW **Dual** Problem

$$\max_{\alpha} \quad -\frac{1}{2} \sum_{c=1}^C \left\| X\alpha_c - \underbrace{\frac{1}{C} \sum_{\tilde{c}} X\alpha_{\tilde{c}}}_{\text{mean}} \right\|^2 + \sum_{c,i:y_i=c} \alpha_i$$

$$\text{s.t. } \alpha_{i,y_i} = 0$$

$$0 \leq \alpha_{i,c} \leq C$$

This is the LLW **Dual** Problem

$$\max_{\alpha} \max_{\bar{w}} - \frac{1}{2} \sum_{c=1}^C \left\| X\alpha_c - \underbrace{\frac{1}{C} \sum_{\tilde{c}} X\alpha_{\tilde{c}}}_{=\bar{w}} \right\|^2 + \sum_{c,i:y_i=c} \alpha_i$$

$$\text{s.t. } \alpha_{i,y_i} = 0$$

$$0 \leq \alpha_{i,c} \leq C$$

This is the LLW **Dual** Problem

$$\begin{aligned} & \max_{\alpha, \bar{w}} \sum_c \overbrace{\left[-\frac{1}{2} \|X\alpha_c - \bar{w}\|^2 + \sum_{i:y_i=c} \alpha_i \right]}^{D_c(\alpha_c, \bar{w})} \\ & \text{s.t. } \alpha_{i,y_i} = 0 \\ & 0 \leq \alpha_{i,c} \leq C \end{aligned}$$

LLW: Proposed Algorithm

Algorithm Simple wrapper algorithm

```
1: function SIMPLESOLVE-LLW( $C, X, Y$ )
2:   while not converged do
3:     for  $c = 1..\mathcal{C}$  do in parallel
4:        $\alpha_c \leftarrow \arg \max_{\tilde{\alpha}_c} D_c(\tilde{\alpha}_c, \bar{w})$ 
5:     end for
6:      $\bar{w} \leftarrow \arg \max_w D(\alpha, w)$ 
7:   end while
8: end function
```

Alber, Zimmert, Dogan, and Kloft (2016):
NIPS submitted

LLW: Proposed Algorithm

Algorithm Simple wrapper algorithm

```
1: function SIMPLESOLVE-LLW( $C, X, Y$ )
2:   while not converged do
3:     for  $c = 1..\mathcal{C}$  do in parallel
4:        $\alpha_c \leftarrow \arg \max_{\tilde{\alpha}_c} D_c(\tilde{\alpha}_c, \bar{w})$ 
5:     end for
6:      $\bar{w} \leftarrow \arg \max_w D(\alpha, w)$ 
7:   end while
8: end function
```

Alber, Zimmert, Dogan, and Kloft (2016):
NIPS submitted rejected ;)

LLW: Proposed Algorithm

Algorithm Simple wrapper algorithm

```
1: function SIMPLESOLVE-LLW( $C, X, Y$ )
2:   while not converged do
3:     for  $c = 1..\mathcal{C}$  do in parallel
4:        $\alpha_c \leftarrow \arg \max_{\tilde{\alpha}_c} D_c(\tilde{\alpha}_c, \bar{w})$ 
5:     end for
6:      $\bar{w} \leftarrow \arg \max_w D(\alpha, w)$ 
7:   end while
8: end function
```

Alber, Zimmert, Dogan, and Kloft (2016):
NIPS submitted rejected ;)
PLoS submitted, arXiv:1611.08480

Ok, fine so far with the LLW SVM.
Now, let's look at the **Weston and Watkins (WW)** SVM.

WW: This is How the **Dual** Problem Looks Like

$$\begin{aligned}
 & \max_{\alpha \in \mathbb{R}^{n \times C}} \quad \overbrace{\sum_{c=1}^C \left[-\frac{1}{2} \| -X\alpha_c \|^2 + \sum_{i: y_i \neq c} \alpha_{i,c} \right]} =: D(\alpha) \\
 & \text{s.t.} \quad \forall i : \alpha_{i, y_i} = - \sum_{c: c \neq y_i} \alpha_{i,c}, \\
 & \quad \forall c \neq y_i : 0 \leq \alpha_{i,c} \leq C
 \end{aligned}$$

WW: This is How the **Dual** Problem Looks Like

$$\begin{aligned}
 & \max_{\alpha \in \mathbb{R}^{n \times \mathcal{C}}} \quad \overbrace{\sum_{c=1}^{\mathcal{C}} \left[-\frac{1}{2} \| -X\alpha_c \|^2 + \sum_{i: y_i \neq c} \alpha_{i,c} \right]} =: D(\alpha) \\
 & \text{s.t.} \quad \forall i : \alpha_{i, y_i} = - \sum_{c: c \neq y_i} \alpha_{i,c}, \\
 & \quad \quad \forall c \neq y_i : 0 \leq \alpha_{i,c} \leq C
 \end{aligned}$$

A common strategy to optimize such a dual problem, is to optimize one coordinate after another (“**dual coordinate ascent**”):

- 1 for $i = 1, \dots, n$
- 2 for $c = 1, \dots, \mathcal{C}$
- 3 $\alpha_{i,c} = \max_{\alpha_{i,c}} D(\alpha)$
- 4 end
- 5 end

This is Now the Story...

We optimize $\alpha_{i,c}$ into gradient direction:

$$\frac{\partial}{\partial \alpha_{i,c}} : 1 - (w_{y_i} - w_c)^T x_i$$

Derivative depends only on **two** weight vectors (not all \mathcal{C} many!).

This is Now the Story...

We optimize $\alpha_{i,c}$ into gradient direction:

$$\frac{\partial}{\partial \alpha_{i,c}} : 1 - (w_{y_i} - w_c)^T x_i$$

Derivative depends only on **two** weight vectors (not all \mathcal{C} many!).

Can we exploit this?

Analogy: Soccer League Schedule

We are given a football league (e.g., Bundesliga) with C many teams.

Before the season, we have to decide on a schedule such that each team plays any other team exactly once.

Furthermore, all teams shall play on every matchday so that in total we need only $C - 1$ matchdays.

Example

Bundesliga has $C = 18$ teams.

⇒ $C - 1 = 17$ matchdays (or twice that many if counting home and away matches)

Analogy: Soccer League Schedule

We are given a football league (e.g., Bundesliga) with C many teams.

Before the season, we have to decide on a schedule such that each team plays any other team exactly once.

Furthermore, all teams shall play on every matchday so that in total we need only $C - 1$ matchdays.

Example

Bundesliga has $C = 18$ teams.

⇒ $C - 1 = 17$ matchdays (or twice that many if counting home and away matches)

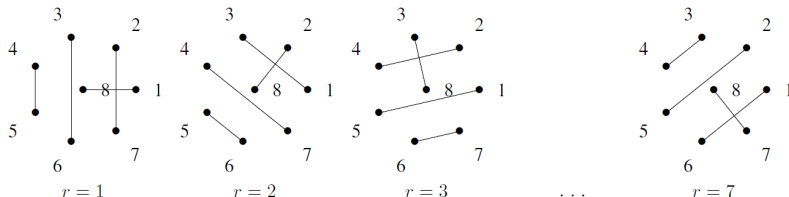
How can we come up with a schedule?

This is a Classical Computer Science Problem...

This is the **1-factorization of a graph** problem.

This is a Classical Computer Science Problem...

This is the **1-factorization of a graph** problem. The solution is known:



Here: $\mathcal{C} = 8$ many teams, 7 matchdays

WW: Proposed Algorithm

Algorithm Simplistic DBCA

wrapper algorithm

```
1: function SIMPLESOLVE-WW( $C, X, Y$ )
2:   while not converged do
3:     for  $r = 1 \dots \mathcal{C} - 1$  do # iterate over "matchdays"
4:       for  $c = 1 \dots \mathcal{C}/2$  do in parallel # iterate over
        "matches"
5:          $(c_i, c_j) \leftarrow$  the two classes ("opposing teams")
6:          $\alpha_{I_{c_i}, c_j}, \alpha_{I_{c_j}, c_i} \leftarrow \arg \max_{\alpha_1, \alpha_2} D_c(\alpha_1, \alpha_2)$ 
7:       end for
8:     end for
9:   end while
10: end function
```

Alber, Zimmert, Dogan, and Kloft (2016): arXiv:1611.08480

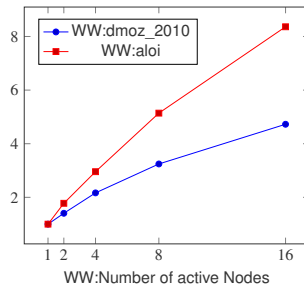
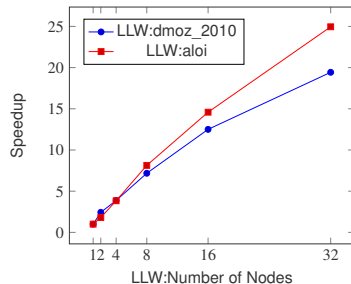
Accuracies

Dataset	# Training	# Test	# Classes	# Features
ALOI	98,200	10,800	1000	128
LSHTCsmall	4,463	1,858	1,139	51,033
DMOZ2010	128,710	34,880	12,294	381,581

Dataset	OVR	CS	WW	LLW
ALOI	0.1824	0.0974	0.0930	0.6560
LSHTCsmall	0.549	0.5919	0.5505	0.9263
DMOZ2010	0.5721	-	0.5432	0.9586

Table: Datasets used in our paper, their properties and best test error over a grid of C values.

Results: Speedup



Open questions

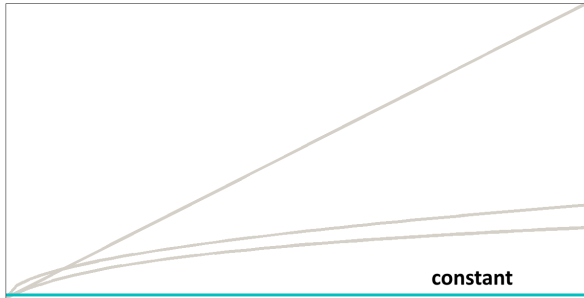
- ▶ higher efficiencies via GPUs?
- ▶ Why does LLW accuracy break?
- ▶ parallelization for CS?

- 1 Introduction
- 2 Distributed Algorithms
- 3 Theory**
- 4 Learning Algorithms
- 5 Conclusion

Theory and Algorithms in Extreme Classification

- ▶ Just saw: **Algorithms** that better handle large number of classes

runtime_algo

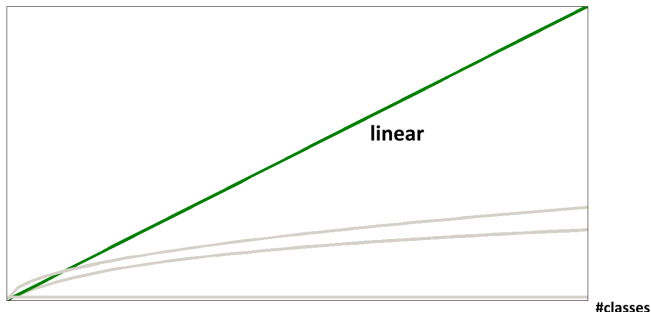


#classes

Theory and Algorithms in Extreme Classification

- ▶ **Theory** not prepared for extreme classification
 - ▶ Data-dependent bounds scale at least **linearly** with the number of classes

(Koltchinskii and Panchenko, 2002; Mohri et al., 2012; Kuznetsov et al., 2014)

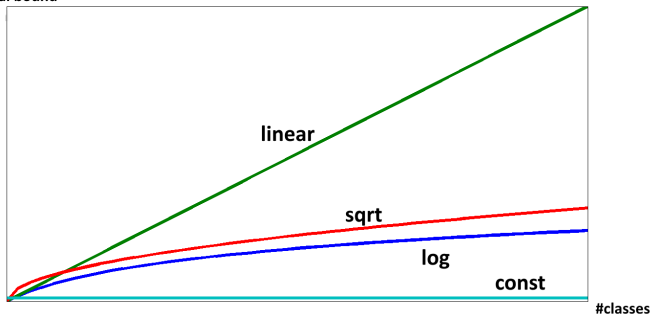


Theory of Extreme Classification

Questions

- ▶ Can we get bounds with **mild** dependence on #classes?
⇒ Novel algorithms?

theoretical bound

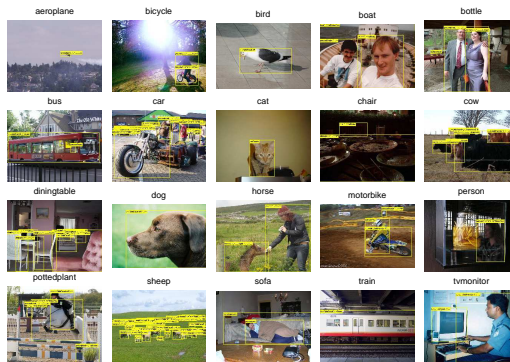


Multi-class Classification

Given:

- ▶ Training data $\underbrace{z_1 = (x_1, y_1), \dots, z_n = (x_n, y_n)}_{\in \mathcal{X} \times \mathcal{Y}} \stackrel{\text{i.i.d.}}{\sim} P$

- ▶ $\mathcal{Y} := \{1, 2, \dots, \mathcal{C}\}$
- ▶ \mathcal{C} = number of classes



Formal Problem Setting

Aim:

- ▶ Define a hypothesis class H of functions $h = (h_1, \dots, h_c)$
- ▶ Find an $h \in H$ that “predicts well” via

$$\hat{y} := \arg \max_{y \in \mathcal{Y}} h_y(x)$$

Multi-class SVMs:

- ▶ $h_y(x) = \langle \mathbf{w}_y, \phi(x) \rangle$
- ▶ Introduce notion of the **(multi-class) margin**

$$\rho_h(x, y) := h_y(x) - \max_{y': y' \neq y} h_{y'}(x)$$

- ▶ the larger the margin, the better

Want: large expected margin $\mathbb{E} \rho_h(X, Y)$.

Types of Generalization bounds for Multi-class Classification

Data-independent bounds

- ▶ based on covering numbers
(Guermeur, 2002; Zhang, 2004a,b; Hill and Doucet, 2007)
- conservative
 - ▶ unable to adapt to data

Data-dependent bounds

- ▶ based on Rademacher complexity
(Koltchinskii and Panchenko, 2002; Mohri et al., 2012; Cortes et al., 2013; Kuznetsov et al., 2014)
- + tighter
 - ▶ able to capture the real data
 - ▶ computable from the data

Def.: Rademacher and Gaussian Complexity

- ▶ Let $\sigma_1, \dots, \sigma_n$ be independent Rademacher variables (taking only values ± 1 , with equal probability).
- ▶ The **Rademacher complexity** (RC) is defined as

$$\mathfrak{R}(H) := \mathbb{E}_{\sigma} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \boxed{\sigma_i} h(z_i) \right]$$

- ▶ Let $g_1, \dots, g_n \sim N(0, 1)$.
- ▶ The **Gaussian complexity** (GC) is defined as

$$\mathfrak{G}(H) = \mathbb{E}_g \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \boxed{g_i} h(z_i) \right]$$

Interpretation: RC and GC reflect the **ability of the hypothesis class to correlate with random noise**.

Def.: Rademacher and Gaussian Complexity

- ▶ Let $\sigma_1, \dots, \sigma_n$ be independent Rademacher variables (taking only values ± 1 , with equal probability).
- ▶ The **Rademacher complexity** (RC) is defined as

$$\mathfrak{R}(H) := \mathbb{E}_{\sigma} \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n \sigma_i h(z_i) \right]$$

- ▶ Let $g_1, \dots, g_n \sim N(0, 1)$.
- ▶ The **Gaussian complexity** (GC) is defined as

$$\mathfrak{G}(H) = \mathbb{E}_g \left[\sup_{h \in H} \frac{1}{n} \sum_{i=1}^n g_i h(z_i) \right]$$

Interpretation: RC and GC reflect the **ability of the hypothesis class to correlate with random noise**.

Theorem ((Ledoux and Talagrand, 1991))

$$\mathfrak{R}(H) \leq \sqrt{\frac{\pi}{2}} \mathfrak{G}(H) \leq 3 \sqrt{\frac{\pi}{2}} \sqrt{\log n} \mathfrak{R}(H).$$

Existing Data-Dependent Analysis

The key step is estimating $\mathfrak{R}(\{\rho_h : h \in H\})$ induced from the **margin operator** ρ_h and class H .

Existing bounds build on the structural result:

$$\mathfrak{R}(\max\{h_1, \dots, h_{\mathcal{C}}\} : h_j \in H_c, c = 1, \dots, \mathcal{C}) \leq \sum_{c=1}^{\mathcal{C}} \mathfrak{R}(H_c) \quad (1)$$

Existing Data-Dependent Analysis

The key step is estimating $\mathfrak{R}(\{\rho_h : h \in H\})$ induced from the **margin operator** ρ_h and class H .

Existing bounds build on the structural result:

$$\mathfrak{R}(\max\{h_1, \dots, h_C\} : h_j \in H_c, c = 1, \dots, C) \leq \sum_{c=1}^C \mathfrak{R}(H_c) \quad (1)$$

Best known dependence on the number of classes:

- ▶ **quadratic** dependence Koltchinskii and Panchenko (2002); Mohri et al. (2012); Cortes et al. (2013)
- ▶ **linear** dependence Kuznetsov et al. (2014)

Existing Data-Dependent Analysis

The key step is estimating $\mathfrak{R}(\{\rho_h : h \in H\})$ induced from the **margin operator** ρ_h and class H .

Existing bounds build on the structural result:

$$\mathfrak{R}(\max\{h_1, \dots, h_C\} : h_j \in H_c, c = 1, \dots, C) \leq \sum_{c=1}^C \mathfrak{R}(H_c) \quad (1)$$

Best known dependence on the number of classes:

- ▶ **quadratic** dependence Koltchinskii and Panchenko (2002); Mohri et al. (2012); Cortes et al. (2013)
- ▶ **linear** dependence Kuznetsov et al. (2014)

Can we do better?

Existing Data-Dependent Analysis

The key step is estimating $\mathfrak{R}(\{\rho_h : h \in H\})$ induced from the **margin operator** ρ_h and class H .

Existing bounds build on the structural result:

$$\mathfrak{R}(\max\{h_1, \dots, h_C\} : h_j \in H_c, c = 1, \dots, C) \leq \sum_{c=1}^C \mathfrak{R}(H_c) \quad (1)$$

Best known dependence on the number of classes:

- ▶ **quadratic** dependence Koltchinskii and Panchenko (2002); Mohri et al. (2012); Cortes et al. (2013)
- ▶ **linear** dependence Kuznetsov et al. (2014)

Can we do better?

The correlation among class-wise components is ignored.

A New Structural Lemma on Gaussian Complexities

We consider Gaussian complexity.

► We show:

$$\mathfrak{G}(\{\max\{h_1, \dots, h_C\} : h = (h_1, \dots, h_C) \in H\}) \leq \boxed{\frac{1}{n} \mathbb{E}_g \sup_{h=(h_1, \dots, h_C) \in H} \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c(x_i)} . \quad (2)$$

A New Structural Lemma on Gaussian Complexities

We consider Gaussian complexity.

► We show:

$$\mathfrak{G}(\{\max\{h_1, \dots, h_C\} : h = (h_1, \dots, h_C) \in H\}) \leq \frac{1}{n} \mathbb{E}_g \sup_{h=(h_1, \dots, h_C) \in H} \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c(x_i). \quad (2)$$

Core idea: **Comparison inequality** on GPs: (Slepian, 1962)

$$\mathfrak{X}_h := \sum_{i=1}^n g_i \max\{h_1(x_i), \dots, h_C(x_i)\}, \mathfrak{Y}_h := \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c(x_i), \forall h \in H.$$

$$\mathbb{E}[(\mathfrak{X}_\theta - \mathfrak{X}_{\bar{\theta}})^2] \leq \mathbb{E}[(\mathfrak{Y}_\theta - \mathfrak{Y}_{\bar{\theta}})^2] \implies \mathbb{E}[\sup_{\theta \in \Theta} \mathfrak{X}_\theta] \leq \mathbb{E}[\sup_{\theta \in \Theta} \mathfrak{Y}_\theta].$$

Eq. (2) preserves the coupling among class-wise components!

Example on Comparison of the Structural Lemma

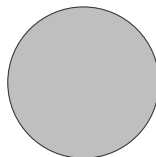
► Consider

$$H := \{(x_1, x_2) \rightarrow (h_1, h_2)(x_1, x_2) = (w_1 x_1, w_2 x_2) : \|(w_1, w_2)\|_2 \leq 1\}$$

► For the function class $\{\max\{h_1, h_2\} : h = (h_1, h_2) \in H\}$,

$$\sup_{(h_1, h_2) \in H} \sum_{i=1}^n \sigma_i h_1(x_i) + \sup_{(h_1, h_2) \in H} \sum_{i=1}^n \sigma_i h_2(x_i)$$

$$\sup_{(h_1, h_2) \in H} \sum_{i=1}^n [g_{i1} h_1(x_i) + g_{i2} h_2(x_i)]$$



Preserving the coupling means supremum in a smaller space!

Estimating Multi-class Gaussian Complexity

- Consider a **vector-valued** function class defined by

$$H := \{h^{\mathbf{w}} = (\langle \mathbf{w}_1, \phi(x) \rangle, \dots, \langle \mathbf{w}_c, \phi(x) \rangle) : f(\mathbf{w}) \leq \Lambda\},$$

where f is **β -strongly convex** w.r.t. $\|\cdot\|$

- $f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) - \frac{\beta}{2}\alpha(1 - \alpha)\|x - y\|^2$.

Theorem

$$\frac{1}{n} \mathbb{E}_{\mathbf{g}} \sup_{h^{\mathbf{w}} \in H} \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c^{\mathbf{w}}(x_i) \leq \frac{1}{n} \sqrt{\frac{2\pi\Lambda}{\beta} \mathbb{E}_{\mathbf{g}} \sum_{i=1}^n \left\| \left(g_{ic} \phi(x_i) \right)_{c=1}^C \right\|_*^2}, \quad (3)$$

where $\|\cdot\|_*$ is the **dual norm** of $\|\cdot\|$.

Features of the complexity bound

- ▶ Applies to a **general** function class defined through a strongly-convex regularizer f
- ▶ Class-wise components h_1, \dots, h_C are correlated through the term $\left\| \left(g_{ic} \phi(x_i) \right)_{c=1}^C \right\|_*^2$
- ▶ Consider class $H_{p,\Lambda} := \{h^{\mathbf{w}} : \|\mathbf{w}\|_{2,p} \leq \Lambda\}$, $(\frac{1}{p} + \frac{1}{p^*} = 1)$; then:

$$\frac{1}{n} \mathbb{E}_{\mathbf{g}} \sup_{h^{\mathbf{w}} \in H_{p,\Lambda}} \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c^{\mathbf{w}}(x_i) \leq \frac{\Lambda}{n} \sqrt{\sum_{i=1}^n k(x_i, x_i)} \times \begin{cases} \sqrt{e} (4 \log C)^{1 + \frac{1}{2 \log C}}, & \text{if } p^* \geq 2 \log C, \\ (2p^*)^{1 + \frac{1}{p^*}} \boxed{C^{\frac{1}{p^*}}}, & \text{otherwise.} \end{cases}$$

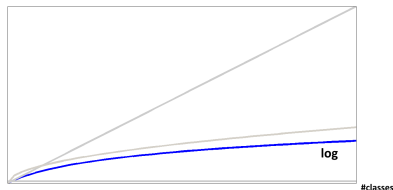
The dependence is **sublinear** for $1 \leq p \leq 2$, and even **logarithmic** when p approaches to 1!

ℓ_p -norm MC-SVM

- Consider class $H_{p,\Lambda} := \{h^{\mathbf{w}} : \|\mathbf{w}\|_{2,p} \leq \Lambda\}$, $(\frac{1}{p} + \frac{1}{p^*} = 1)$; then:

$$\frac{1}{n} \mathbb{E}_g \sup_{h^{\mathbf{w}} \in H_{p,\Lambda}} \sum_{i=1}^n \sum_{c=1}^C g_{ic} h_c^{\mathbf{w}}(x_i) \leq \frac{\Lambda}{n} \sqrt{\sum_{i=1}^n k(x_i, x_i)} \times \begin{cases} \sqrt{e} (4 \log C)^{1 + \frac{1}{2 \log C}}, & \text{if } p^* \geq 2 \log C, \\ (2p^*)^{1 + \frac{1}{p^*}} \boxed{C^{\frac{1}{p^*}}}, & \text{otherwise.} \end{cases}$$

The dependence is **sublinear** for $1 \leq p \leq 2$, and even **logarithmic** when p approaches to 1!



Future Directions

Theory: A data-dependent bound **independent** of the class size?

Future Directions

Theory: A data-dependent bound **independent** of the class size?

- ⇒ Need more powerful structural result on Gaussian complexity for functions induced by **maximum operator**.
- ▶ Might be worth to look into **ℓ_∞ -norm covering numbers**.

Reference: Lei, Dogan, Binder, and Kloft (NIPS 2015);
Journal submission forthcoming

- 1 Introduction
- 2 Distributed Algorithms
- 3 Theory
- 4 Learning Algorithms**
- 5 Conclusion

ℓ_p -norm Multi-class SVM

Motivated by the **mild dependence** on \mathcal{C} as $p \rightarrow 1$, we consider

$(\ell_p$ -norm) Multi-class SVM, $1 \leq p \leq 2$

$$\begin{aligned} \min_{\mathbf{w}} \quad & \frac{1}{2} \left[\sum_{c=1}^{\mathcal{C}} \|\mathbf{w}_c\|_2^p \right]^{\frac{2}{p}} + C \sum_{i=1}^n (1 - t_i)_+, \\ \text{s.t.} \quad & t_i = \langle \mathbf{w}_{y_i}, \phi(x_i) \rangle - \max_{y: y \neq y_i} \langle \mathbf{w}_y, \phi(x_i) \rangle, \end{aligned} \tag{P}$$

Empirical Results

Empirical Results:


Method / Dataset	Sector	News 20	Rcv1	Birds 50	Caltech 256
ℓ_p -norm MC-SVM	94.2 \pm 0.3	86.2 \pm 0.1	85.7 \pm 0.7	27.9 \pm 0.2	56.0 \pm 1.2
Crammer & Singer	93.9 \pm 0.3	85.1 \pm 0.3	85.2 \pm 0.3	26.3 \pm 0.3	55.0 \pm 1.1

Proposed ℓ_p -norm MC-SVM consistently better on benchmark datasets.

Wait... I performed this Experiment:

Wait... I performed this Experiment:

- So I took the DMOZ2010 dataset
(Aim: categorize new webpages)

 open directory project

In partnership with
Aol Search.

[about dmoz](#) | [dmoz blog](#) | [suggest URL](#) | [help](#) | [link](#) | [editor login](#)

Search

advanced

[Arts](#)
Movies, Television, Music...

[Games](#)
Video Games, RPGs, Gambling...

[Kids and Teens](#)
Arts, School Time, Teen Life...

[Reference](#)
Maps, Education, Libraries...

[Shopping](#)
Clothing, Food, Gifts...

[World](#)
Català, Dansk, Deutsch, Español, Français, Italiano, 日本語, Nederlands, Polski, Русский, Svenska...

[Business](#)
Jobs, Real Estate, Investing...

[Health](#)
Fitness, Medicine, Alternative...

[News](#)
Media, Newspapers, Weather...

[Regional](#)
US, Canada, UK, Europe...

[Society](#)
People, Religion, Issues...

[Computers](#)
Internet, Software, Hardware...


[Home](#)
Family, Consumers, Cooking...

[Recreation](#)
Travel, Food, Outdoors, Humor...

[Science](#)
Biology, Psychology, Physics...

[Sports](#)
Baseball, Soccer, Basketball...

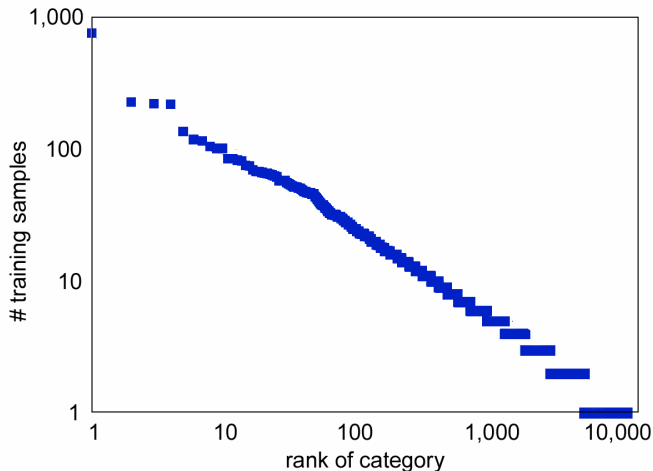
[Become an Editor](#) | Help build the largest human-edited directory of the web



Copyright © 2013 Netscape

Wait... I performed this Experiment:

- ▶ OVR-SVM, Train=128,710, Test=34,880; Result:



27% of classes never used in prediction

New Learning Algorithm

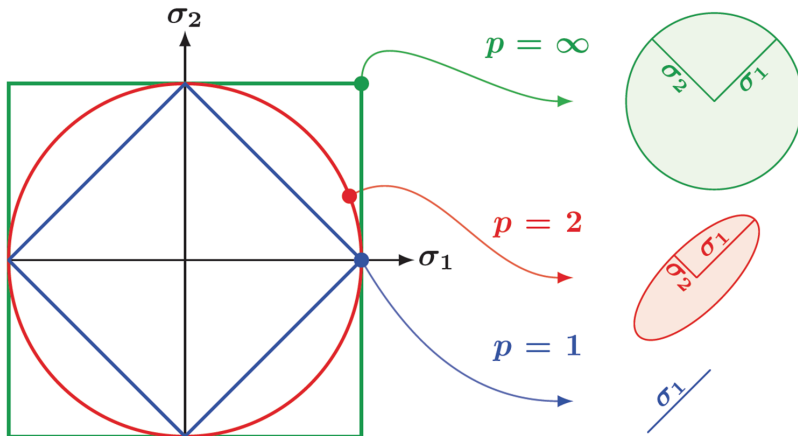
Schatten-SVM

$$\min_{W=(w_1, \dots, w_C)} \frac{1}{2} \sum_c \underbrace{\|W\|_{S_p}^2}_{\text{Schatten norm}} + C \sum_{i=1}^n \left[\max_{c \neq y_i} l((w_{y_i} - w_c)^T x_i) \right]$$

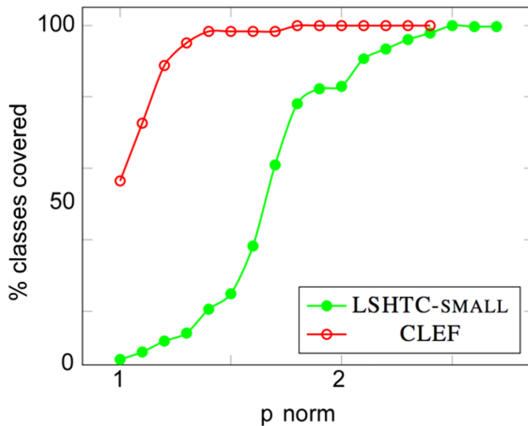
Schatten- p norm

$$\|W\|_{S_p} := \sqrt[p]{\sum_i \sigma_i^p(\sqrt{W^T W})}$$

Geometry of Schatten Norm



Schatten-norm Parameter p Controls coverage



Results

Dataset		Schatten-SVM	OvR	CS-SVM	HR-SVM	HR-LR	TD-SVM
CLEF	Macro-F1	58.42 (52.20)	53.11	57.17	53.92	55.83	32.32
	Micro-F1	80.21 (78.82)	78.92	79.94	80.02	80.12	70.11
	Coverage	90.48 (85.71)	87.30	88.93			
LSHTC-SMALL	Macro-F1	30.10 (30.12)	26.89	28.22	28.94	28.12	20.01
	Micro-F1	46.12 (45.85)	43.34	45.77	45.31	44.94	38.48
	Coverage	60.66 (61.54)	54.52	55.87			
WIKI-2011	Macro-F1	30.29	25.13	27.35	-	-	-
	Micro-F1	44.86	39.07	43.47	-	-	-
	Coverage	74.58	61.51	67.90			
DMOZ-2010	Macro-F1	32.04	31.27	32.64	33.12	32.42	22.30
	Micro-F1	44.12	45.12	45.36	46.02	45.84	38.64
	Coverage	68.57	63.82	64.50			

Future Directions

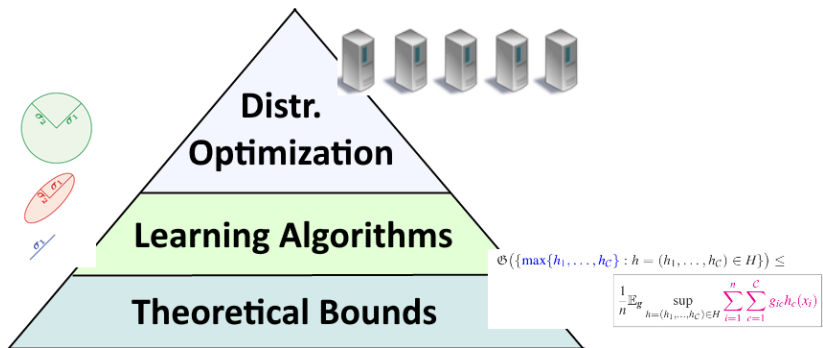
Algorithms: New models & efficient solvers

- ▶ **Novel models** motivated by theory
 - ▶ top-k MC-SVM (Lapin et al., 2015)
- ▶ Analyze $p > 2$ regime
- ▶ Extensions to **multi-label** learning

- 1 Introduction
- 2 Distributed Algorithms
- 3 Theory
- 4 Learning Algorithms
- 5 Conclusion**

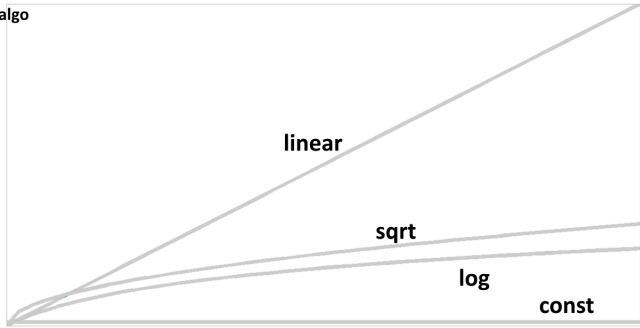
Conclusion

Extreme Classification



Conclusion

theoretical bound &
runtime_algo



#classes

Refs I

- M. Alber, J. Zimmert, U. Dogan, and M. Kloft. Distributed Optimization of Multi-Class SVMs. **in submission**, 2016.
- C. Cortes, M. Mohri, and A. Rostamizadeh. Multi-class classification with maximum margin multiple kernel. In **ICML-13**, pages 46–54, 2013.
- K. Crammer and Y. Singer. On the algorithmic implementation of multiclass kernel-based vector machines. **Journal of Machine Learning Research**, 2:265–292, 2002.
- Y. Guermeur. Combining discriminant models with new multi-class svms. **Pattern Analysis & Applications**, 5(2): 168–179, 2002.
- S. I. Hill and A. Doucet. A framework for kernel-based multi-category classification. **Journal of Artificial Intelligence Research**, 30(1):525–564, 2007.
- V. Koltchinskii and D. Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. **Annals of Statistics**, pages 1–50, 2002.
- V. Kuznetsov, M. Mohri, and U. Syed. Multi-class deep boosting. In **Advances in Neural Information Processing Systems**, pages 2501–2509, 2014.
- M. Lapin, M. Hein, and B. Schiele. Top-k multiclass SVM. **CoRR**, abs/1511.06683, 2015. URL <http://arxiv.org/abs/1511.06683>.
- M. Ledoux and M. Talagrand. **Probability in Banach Spaces: isoperimetry and processes**, volume 23. Springer, Berlin, 1991.
- Y. Lee, Y. Lin, and G. Wahba. Multicategory support vector machines: Theory and application to the classification of microarray data and satellite radiance data. **Journal of the American Statistical Association**, 99(465):67–82, 2004.
- Y. Lei, Ü. Dogan, A. Binder, and M. Kloft. Multi-class svms: From tighter data-dependent generalization bounds to novel algorithms. In **Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada**, pages 2035–2043, 2015. URL <http://papers.nips.cc/paper/6012-multi-class-svms-from-tighter-data-dependent-generalization-bounds-to-novel-algorithms>
- M. Mohri, A. Rostamizadeh, and A. Talwalkar. **Foundations of machine learning**. MIT press, 2012.

Refs II

- D. Slepian. The one-sided barrier problem for gaussian noise. **Bell System Technical Journal**, 41(2):463–501, 1962.
- J. Weston and C. Watkins. Support vector machines for multi-class pattern recognition. In M. Verleysen, editor, **Proceedings of the Seventh European Symposium On Artificial Neural Networks (ESANN)**, pages 219–224. Evere, Belgium: d-side publications, 1999.
- T. Zhang. Class-size independent generalization analysis of some discriminative multi-category classification. In **Advances in Neural Information Processing Systems**, pages 1625–1632, 2004a.
- T. Zhang. Statistical analysis of some multi-category large margin classification methods. **The Journal of Machine Learning Research**, 5:1225–1251, 2004b.