Fast optimization of Multithreshold Entropy Linear Classifier

Rafal Jozefowicz, Wojciech Czarnecki

19 February 2015

Rafal Jozefowicz, Wojciech Czarnecki

19 February 2015 1 / 22

Idea

Classification procedure:

- Find the best projection w for the training set
- Perform kernel density estimation for each class on the projected data
- Classify samples based on the KDE values

I will focus on the first part of the problem as the other two are not as computationally expensive.



Find the best projection w for the training set

$$\begin{split} \mathsf{maximize}_{w \in \mathbb{R}^d} \ \mathsf{D}_{CS}(\llbracket w^\mathsf{T} X_- \rrbracket, \llbracket w^\mathsf{T} X_+ \rrbracket) \\ \mathsf{subject to} \ \lVert w \rVert = 1 \end{split}$$

where [X] is the kernel density estimation of X.

Find the best projection w for the training set

$$\begin{split} \mathsf{maximize}_{w \in \mathbb{R}^d} \ \mathsf{D}_{CS}(\llbracket w^{\mathsf{T}} X_{-} \rrbracket, \llbracket w^{\mathsf{T}} X_{+} \rrbracket) \\ \mathsf{subject to} \ \lVert w \rVert = 1 \end{split}$$

where [X] is the kernel density estimation of X.

$$D_{CS}(f,g) = -2\log \int fg(x)dx + \log \int f^{2}(x)dx + \log \int g^{2}(x)dx$$
$$= -2\log ip^{\times}(f,g) + \log ip^{\times}(f,f) + \log ip^{\times}(g,g)$$
$$ip^{\times}(w) = \frac{1}{\sqrt{2\pi V(w)|X||Y|}} \sum_{x,y} \exp\left(-\frac{\langle w, x-y\rangle^{2}}{2V(w)}\right)$$
$$V(w) = \gamma^{2}(\frac{4}{3})^{\frac{1}{25}}(|X|^{-\frac{7}{5}}\sum_{x}\langle x-\hat{X},w\rangle^{2} + |Y|^{-\frac{7}{5}}\sum_{y}\langle y-\hat{Y},w\rangle^{2})$$

where
$$\mathsf{ip}^{ imes}(w) = \mathsf{ip}^{ imes}(\llbracket w^{\mathsf{T}}X_{-} \rrbracket, \llbracket w^{\mathsf{T}}X_{+} \rrbracket).$$

- 31

Computational complexity

Computing $ip^{\times}(w)$ is computationally expensive as it requires O(|X||Y|) steps.

・ロン ・四 ・ ・ ヨン ・ ヨン

Approximations!

We investigated two approximations of $ip^{\times}(\cdot)$:

- Sorting and Discarding ignores pairs of points that are too far away from each other as they don't contribute much to the final value. The maximum distance is determined dynamically for each *w* by the approximation factor.
- Binning performs binning of the projected points, so that those located near each other are approximated by their empirical mean. The number of bins is determined dynamically for each w by the approximation factor.

Both of them are parameterized by ϵ , which bounds the absolute error made by the algorithms vs ip[×](·).

Sorting and Discarding

$$\mathsf{ip}_{\mathsf{sort}}^{\times}(w,\epsilon) = \frac{1}{\sqrt{2\pi V(w)|X||Y|}} \sum_{x,y \,:\, |\mathbf{x}-\mathbf{y}| < \mathsf{width}(\epsilon)} \exp\left(-\frac{\langle w, x-y \rangle^2}{2V(w)}\right)$$

Theorem

Using adaptive sorting and discarding with distance threshold in each iteration of at least

$$\sqrt{\max\left\{0,-V(w)\ln\left(2\epsilon^2\pi V(w)
ight)
ight\}}$$

leads to the computation of the ip^{\times} function with at most ϵ error.

Sorting and Discarding

$$\mathsf{ip}_{\mathsf{sort}}^{\times}(w,\epsilon) = \frac{1}{\sqrt{2\pi V(w)|X||Y|}} \sum_{x,y \,:\, |\mathbf{x}-\mathbf{y}| < \mathsf{width}(\epsilon)} \exp\left(-\frac{\langle w, x-y\rangle^2}{2V(w)}\right)$$

The algorithm works best for sparse projections. It can be easily implemented in $O(n \log n + |\{x, y : |x - y| < width(\epsilon)\}|)$ time and linear memory.

Binning

$$\mathsf{ip}_{\mathsf{bin}}^{\times}(w,\epsilon) = \frac{1}{\sqrt{2\pi V(w)|X||Y|}} \sum_{x,y} \exp\left(-\frac{(\langle w,x\rangle_b - \langle w,y\rangle_b)^2}{2V(w)}\right)$$

Theorem

Using adaptive binning technique with bin width in each iteration at most

$$\sqrt{-2V(w)\ln\left(\max\left\{0,1-\epsilon\sqrt{2\pi V(w)}
ight\}
ight)}$$

leads to the computation of the ip^{\times} function with at most ϵ error.

Binning

$$\mathsf{ip}_{\mathsf{bin}}^{\times}(w,\epsilon) = \frac{1}{\sqrt{2\pi V(w)|X||Y|}} \sum_{x,y} \exp\left(-\frac{(\langle w, x \rangle_b - \langle w, y \rangle_b)^2}{2V(w)}\right)$$

The algorithm works best for dense projections. The time complexity is quadratic in the number of bins.

3

・ロン ・四 ・ ・ ヨン ・ ヨン



◆□> ◆圖> ◆臣> ◆臣> □臣

Below are plots of the discarding threshold and bin width as the functions of the acceptable error ϵ .



э

-

• • • • •

Optimization on the sphere

For better numerical stability, the optimization should be performed on the sphere. By adding a custom regularization term we "guide" the optimization method to stay close to it.

Theorem

Given arbitrary sets $X_-, X_+ \subset \mathbb{R}^d$ and corresponding $D_{CS}(w) = D_{CS}(\llbracket w^T X_- \rrbracket, \llbracket w^T X_+ \rrbracket)$ function we have:

$$d := \max_{\|w\|=1} D_{CS}(w) = \max_{w} D_{CS}(w) - (\|w\|^2 - 1)^2$$

and

$$\{w: \|w\| = 1 \land D_{CS}(w) = d\} = \{w: D_{CS}(w) - (\|w\|^2 - 1)^2 = d\}.$$

イロト 不得下 イヨト イヨト

Optimization on the sphere

The new objective function allows us to use optimization techniques that are not designed to work on the sphere. In particular, we used L-BFGS and Conjugate Gradients. We don't need any additional hyperparameters to be fitted.

We evaluated proposed approximations on 10 datasets from UCI and libSVM's repositories. Both D_{CS} and its approximations are coded in Python using numpy and scipy.

Mean ratio of exp calls between approximated technique and original method during optimizations.

method	C	G	L-BFGS-B					
name	bin	dist	bin	dist				
australian	0.11	0.44	0.11	0.45				
breast-cancer	0.10	0.46	0.10	0.46				
diabetes	0.21	0.56	0.22	0.54				
fourclass	0.19	0.51	0.19	0.49				
german.numer	0.15	0.47	0.19	0.46				
heart	0.29	0.47	0.26	0.47				
ionosphere	0.25	0.55	0.24	0.54				
liver-disorders	0.29	0.65	0.31	0.67				
sonar	0.32	0.53	0.29	0.50				
splice	0.19	0.44	0.16	0.43				

Sorting and Discarding - Results

Comparison of the cross validation BAC scores between given γ hyperparameter of D_{CS} (x-axis), accepted error ϵ (y-axis). Positive values (and corresponding red colors) represent decrease in BAC score while negative values and corresponding blue colors – increase after using approximated method.

	australian						breast-cancer					diabetes					fourclass						german				
4.5	1.44%	0.90%	-0.88%	-1.41%	-9.18% +5	0.99%	-1.54%	-2.33%	-3.22%	-0.67% +5	0.38%	-2.67%	0.80%	-3.68%	-4.79%	-6.55%	-5.40%	-0.56%	-2.15%	-4.06% **	0.53%	-1.42%	-6.43%	-9.33%	-19.50%		
4.2	0.43%	-1.78%	0.80%	-0.07%	-0.08% **	1.17%	0.19%	-0.25%	-2.44%	-1.47% **	-2.25%	-2.83%	-1.58%	-4.05%	-3.28%	-0.57%	-2.78%	0.13%	-0.82%	-0.58% 12	1.23%	-2.51%	-6.26%	-6.02%	-2.97%		
6.1	-0.36%	0.27%	0.14%	-0.04%	-0.17% ***	1.19%	-0.41%	-0.49%	-0.89%	-0.57% 41	1.40%	-1.50%	-0.61%	-2.97%	-0.55% -	-0.78%	-4.84%	0.16%	0.61%	0.28%	4.39%	-1.26%	-2.41%	-5.79%	-1.89%		
0.05	0.57%	0.85%	-3.56%	-0.18%	-0.12% ess	1.30%	0.59%	-0.55%	-2.23%	-0.22% ess	-0.38%	-0.68%	0.99%	-2.70%	-1.32% •	-0.70%	-2.32%	0.01%	0.84%	-1.08% ess	3.38%	-3.72%	-1.59%	-4.27%	-4.00%		
0.03	-2.15%	-0.52%	-1.41%	-0.14%	-0.34% ***	1.07%	0.24%	-0.02%	-0.27%	-0.05% ***	0.94%	0.70%	0.82%	0.10%	0.17% at	-0.72%	0.00%	-0.71%	-0.14%	0.20% ***	0.60%	-4.95%	-4.59%	-1.82%	-3.99%		
0.02	0.22%	0.34%	-1.42%	-0.01%	-0.12% ***	0.90%	0.27%	-0.17%	-2.87%	-0.05% ***	0.16%	-1.45%	-0.10%	0.17%	-0.88% •	-0.66%	0.03%	-0.12%	-0.01%	0.74% 012	0.96%	-7.07%	-3.57%	-3.62%	-2.45%		
6.01	1.01%	-1.30%	0.29%	-0.18%	-0.21% ***	0.99%	-0.31%	-0.56%	-0.22%	-0.06% ***	0.72%	-0.29%	3.07%	-0.21%	-0.13% **	-0.81%	-0.03%	0.78%	-0.65%	0.10% and	2.89%	-3.48%	-3.74%	-2.74%	-3.15%		
	0.1	0.5	heart	15	2.0	0.1	os ic	nospher	15 'e	2.0	0.1	live	er-disord	ers	2.0	0.1	0.5	sonar	15	2.0	0.1	0.5	splice	15	2.0		
4.5	-3.06%	-5.20%	-2.58%	-2.83%	-13.47% +3	-1.24%	-8.93%	-8.08%	-13.91%	-25.65% «>	-3.01%	-6.04%	-4.36%	-9.16%	-9.45%	-0.51%	0.13%	-3.98%	-12.34%	-14.37% as	10.61%	-4.12%	-8.54%	-6.85%	-22.07%		
4.2	1.85%	0.48%	-2.25%	1.55%	-0.15% ===	0.18%	-4.35%	-3.04%	-4.73%	-4.52% +2	-2.55%	-5.29%	-5.97%	-4.53%	-2.93%	2 -3.06%	5.35%	-2.10%	-17.60%	-15.46% =>	8.28%	-2.57%	-4.67%	0.21%	-4.17%		
6.1	1.93%	-0.03%	-1.24%	-2.40%	0.48% ***	2.31%	-3.91%	-3.78%	-1.00%	-1.71% 41	-1.73%	-2.71%	-2.92%	-0.20%	-2.97%	-3.04%	-2.85%	-4.28%	-0.82%	-4.33% 41	3.32%	-2.36%	-4.03%	0.32%	-0.38%		
0.05	0.13%	-2.55%	-0.40%	-3.26%	-0.30% ess	-0.82%	-4.86%	-2.74%	-1.05%	2.40% •==	-0.77%	-2.15%	-1.81%	-0.82%	1.25% 0	-2.97%	13.66%	-4.84%	7.44%	-0.59% ess	3.38%	-2.32%	-2.59%	-0.54%	-0.58%		
0.03	0.35%	-0.94%	0.15%	-0.05%	-0.91% ***	3.50%	-2.26%	-1.28%	1.44%	0.90% ***	-2.23%	-4.17%	-6.24%	-0.42%	1.62% **	-0.85%	7.75%	-2.10%	5.99%	-0.03% ==>	7.61%	2.80%	-0.66%	-0.20%	0.08%		
6.02	-1.32%	-0.92%	-0.58%	0.28%	-2.54% ***	1.05%	-2.90%	-2.76%	0.51%	3.13% 082	1.74%	-4.89%	-1.73%	0.98%	0.79% •	-3.25%	11.57%	0.87%	-4.96%	-7.63% ***	3.24%	-1.51%	-0.60%	-0.56%	-0.67%		
6.85	1.40%	-0.74%	0.20%	0.31%	-1.09% 685	2.26%	-6.23%	-0.72%	1.01%	1.89% •**	-3.39%	-1.03%	-1.16%	-4.98%	-0.16% =	-4.75%	8.99%	-0.05%	-5.10%	2.99% •**	7.48%	2.25%	-0.34%	-0.25%	-1.97%		
	0.1	0.5	1.0	15	2.0	0.1	0.5	1.0	1.5	2.0	0.1	0.5	1.0	15	2.0	0.1	0.5	1.0	15	2.0	0.1	0.5	10	1.5	2.0		

Binning - Results

Comparison of the cross validation BAC scores between given γ hyperparameter of D_{CS} (x-axis), accepted error ϵ (y-axis). Positive values (and corresponding red colors) represent decrease in BAC score while negative values and corresponding blue colors – increase after using approximated method.

	australian						breast-cancer					diabetes					fourclass						german				
0.5	-1.27%	2.12%	0.48%	-0.35%	-0.26% +5	1.64%	1.24%	0.32%	-1.29%	-3.56% +3	1.40%	-0.85%	-0.57%	-4.17%	-11.94%	85	-7.50%	-4.74%	1.78%	-3.15%	1.70% .:	6.85%	-6.15%	-7.19%	-9.41%	-18.10%	
0.2	2.76%	2.00%	0.39%	-0.35%	-0.50% - 12	0.91%	0.10%	0.18%	-0.35%	-0.25% ===	3.57%	-0.74%	0.55%	-2.75%	-5.33%	8.2	-3.25%	-2.93%	2.19%	1.84%	2.03% .	4.97%	-1.97%	-5.03%	-5.82%	-5.21%	
6.1	2.63%	0.92%	0.63%	-0.17%	-0.03% - ==	1.07%	0.79%	0.01%	-0.41%	-0.39% **	1.24%	-0.84%	0.72%	-3.04%	-5.25%	63	-4.29%	-4.57%	3.33%	0.25%	-0.34% ::	6.63%	-8.87%	-7.34%	-7.01%	-8.66%	
0.05	2.32%	1.66%	0.39%	-0.20%	-0.26% ess	1.38%	0.54%	0.02%	-0.32%	-0.22% ess	-2.52%	-0.74%	0.94%	-2.50%	-2.73%	0.05	-4.47%	-4.03%	2.85%	1.35%	0.67% as	4.78%	-2.27%	-4.29%	-6.01%	-6.09%	
0.03	1.97%	1.49%	1.05%	-0.35%	-0.26% ***	1.25%	0.37%	0.02%	-0.24%	-0.17% 683	-0.16%	-0.91%	-0.50%	-2.55%	-2.86%	6.83	-6.13%	-6.59%	2.68%	0.25%	0.71% es	6.34%	-4.89%	-4.27%	-5.97%	-4.81%	
0.02	1.97%	1.90%	0.72%	-0.35%	-0.08% 012	1.06%	0.78%	0.10%	-0.35%	-0.30% 612	2.54%	-0.82%	0.86%	-2.63%	-2.98%	0.02	-6.59%	-5.45%	2.79%	1.41%	1.73%	4.99%	-1.36%	-7.23%	10.45%	-6.47%	
0.01	1.15%	1.53%	0.51%	-0.35%	-0.26% ***	0.54%	0.90%	0.26%	-0.30%	-0.22% 681	0.56%	-0.74%	0.87%	-3.47%	-2.84%		10.70%	-4.47%	2.36%	1.94%	0.16% est	7.29%	-5.20%	-7.00%	-6.84%	-6.72%	
	0.1	0.5	heart	15	2.0	0.1	0.5 ic	nosphe	re ¹⁵	2.0	0.1	0.5 live	r-disord	ers	2.0		0.1	0.5	sonar	15	2.0	0.1	0.5	splice	15	2.0	
4.5	1.07%	2.37%	1.59%	-4.87%	-8.03% +5	1.53%	-1.39%	-9.55%	-9.82%	-3.29% +3	-1.81%	-2.66%	-0.85%	-4.69%	-4.75%	-	-7.24%	1.53%	-3.04%	2.64%	-6.53% +	6.29%	4.58%	2.11%	-6.20%	-8.82%	
0.2	-0.08%	1.76%	2.15%	1.55%	-1.75% - 1.2	-2.29%	-2.58%	-2.13%	-6.56%	-5.98% 12	4.09%	-0.40%	-3.26%	-6.61%	1.28%	6.2	-2.06%	-0.56%	3.13%	1.14%	3.30%	5.04%	4.62%	1.87%	1.25%	1.05%	
6.1	1.99%	2.22%	1.80%	1.10%	-0.43% - 11	-1.28%	-1.07%	-1.13%	-3.08%	-3.71% ==	-1.26%	-1.34%	-1.32%	0.32%	-1.11%	6.1	-6.20%	1.16%	4.71%	5.32%	-0.09% **	4.34%	4.64%	2.19%	1.08%	1.27%	
0.05	2.90%	1.83%	2.12%	0.97%	-1.35% ess	-0.76%	-2.29%	-2.07%	-4.34%	-6.74% -15	-7.74%	-1.97%	-1.77%	-5.55%	-3.25%	0.05	3.36%	0.50%	0.21%	3.97%	1.12% 48	4.36%	4.52%	1.90%	1.31%	0.86%	
0.03	2.55%	1.66%	1.92%	0.99%	-0.18% ***	-3.34%	-2.65%	0.90%	-4.26%	-5.22% •**	-6.14%	-0.94%	-1.68%	0.96%	-1.77% (6.83	-6.02%	0.40%	3.23%	3.37%	1.33% •==	3.74%	4.20%	1.79%	1.21%	1.24%	
0.02	0.78%	1.83%	1.47%	0.97%	-0.20% 0.12	-1.84%	-2.24%	-2.25%	-5.24%	-4.86% 412	4.02%	-1.00%	-1.17%	0.98%	-3.31%	6.82	-8.00%	-3.70%	5.73%	4.94%	0.81% es	5.49%	4.09%	2.09%	1.31%	1.15%	
0.01	1.50%	0.94%	1.29%	0.76%	0.18% ***	-3.09%	-3.31%	-1.41%	-2.94%	-4.72% •==	1.17%	-2.37%	-1.11%	-3.49%	0.65%	6.81	0.93%	-2.06%	5.67%	3.68%	0.74% as	4.64%	4.47%	1.86%	1.34%	1.21%	
	0.1	0.5	1.0	1.5	2.0	0.1	0.5	1.0	1.5	2.0	0.1	0.5	1.0	1.5	2.0		0.1	0.5	1.0	1.5	2.0	0.1	0.5	1.0	15	2.0	

Number of optimization methods' iterations.

method		CG		L-BFGS-B				
name	bin	D _{CS}	dist	bin	D _{CS}	dist		
australian	4	36	22	11	39	37		
breast-cancer	4	35	8	6	39	14		
diabetes	3	30	20	18	36	29		
fourclass	4	12	10	6	15	14		
german.numer	7	60	32	7	58	38		
heart	3	40	19	12	34	20		
ionosphere	5	600	216	18	384	152		
liver-disorders	4	30	22	22	43	30		
sonar	4	262	115	15	139	100		
splice	4	92	26	14	65	41		

(日) (四) (三) (三) (三)

Number of optimization methods' iterations.

It's important to note that the number of iterations is not equal to the total number of evaluations of the function. It suggests, though, that when using approximations, the problem becomes simpler than the baseline.

イロト 不得下 イヨト イヨト

Conclusions

- We proposed two simple approximation schemes for faster computation of MELC objective function and its gradient.
- We showed how to efficiently use existing off-the-shelf optimizers by a simple change of the objective function while at the same time still work near the unit sphere. Without the regularization term, the norm of w tends to explode on some datasets and making it numerically unstable.
- Both algorithms significantly reduce the mean number of exp calls while not sacrificing the resulting accuracy.
- Our experiments suggest that the approximations act like some kind of regularization of the classifier.

Combining methods

The algorithms complement each other and it is easy to determine, which will be faster for a given projection. This can further improve experimental results and make it more applicable for larger datasets.