To Search or to Filter? Consistent Label Tree Classifiers for Efficient Classification

Arkadiusz Jachnik

Intelligent Decision Support Systems Laboratory (IDSS) Poznań University of Technology, Poland



PL-SIGML, Wroclaw, March 17-18, 2014

This is a join work with Krzysztof Dembczynski.

- Multi-class classification (MCC): a classification task with one target variable belonging to more than two classes.
- Multi-label classification (MLC): variant of the classification problem where none, one or more class labels are assigned to single instances simultaneously.

Example

Multi-class classification

Target: forest OR path OR dog OR ...

Multi-label classification

Target 1:	forest	yes/no
Target 2:	path	yes/no
Target 3:	dog	yes/no



• Multi-class classification: For a feature vector x predict accurately one response y using a function h(x):

$$\boldsymbol{x} = (x_1, x_2, \dots, x_p) \xrightarrow{h(\boldsymbol{x})} y \in \{1, \dots, k\}$$

• Multi-label classification: For a feature vector \boldsymbol{x} predict accurately a vector of responses \boldsymbol{y} using a function $\boldsymbol{h}(\boldsymbol{x}) = (h_1(\boldsymbol{x}), h_2(\boldsymbol{x}), \dots, h_m(\boldsymbol{x}))$:

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

Binary relevance: Decomposes the problem to m independent binary classification problems:

Binary relevance: Decomposes the problem to m independent binary classification problems:

$$(\boldsymbol{x}, \boldsymbol{y}) \longrightarrow (\boldsymbol{x}, y = y_i), \quad i = 1, \dots, m$$

$$y_1 \qquad y_2 \qquad \cdots \qquad y_m$$

$$y_1 = 0 \qquad \boldsymbol{x} \qquad \boldsymbol{y}_1 = 1 \qquad \boldsymbol{y}_2 = 0 \qquad \boldsymbol{y}_2 = 1 \qquad \cdots \qquad \boldsymbol{y}_m = 0 \qquad \boldsymbol{y}_m = 1$$

$$(P(y_1 = 0 | \boldsymbol{x})) \qquad (P(y_1 = 1 | \boldsymbol{x})) \qquad (P(y_2 = 0 | \boldsymbol{x})) \qquad (P(y_2 = 1 | \boldsymbol{x})) \qquad (P(y_m = 0 | \boldsymbol{x})) \qquad (P(y_m = 1 | \boldsymbol{x}))$$

Challenges:

- How to model dependencies between labels?
- What about loss function defined over the binary vectors?

Label Tree Classifiers for MLC

- Example: m = 2, $\mathcal{Y} = \{0, 1\}^2$
- $2^m 1$ classifiers on the tree (*m* levels); *m*-classifier trick with *m* classifiers (one per level)
- 2^m leaves of the all possible vectors of responses



Label Tree Classifiers for MCC

- We assign each class an integer from 0 to k-1 and code it by its binary representation on m bits.
- Example: k = 4, $\mathcal{Y} = \{0, 1, 2, 3\}$
- k-1 classifiers on the tree; trick with $\log k$ classifiers (one per level)
- k leaves, one for each class



• The prediction accuracy of *h* is measured in terms of its **risk**, that is, its expected loss

$$L(\boldsymbol{h}, P) = \mathbb{E}\left[\ell(\boldsymbol{Y}, \boldsymbol{h}(\boldsymbol{X}))\right] = \int \ell(\boldsymbol{y}, \boldsymbol{h}(\boldsymbol{x})) \, dP(\boldsymbol{x}, \boldsymbol{y}) \, ,$$

where ℓ is a loss function.

• Here, we consider the 0/1 loss:

$$\ell_{0/1}(\boldsymbol{y},\boldsymbol{h}(\boldsymbol{x})) = [\![\boldsymbol{y} \neq \boldsymbol{h}(\boldsymbol{x})]\!]$$

• The optimal Bayes classifier minimizes the risk:

$$\boldsymbol{h}^*(\boldsymbol{x}) = \operatorname*{arg\,min}_{\boldsymbol{h}\in\mathcal{Y}} \int \sum_{\boldsymbol{y}\in\mathcal{Y}} \ell_{0/1}(\boldsymbol{y},\boldsymbol{h}) P(\boldsymbol{y}\,|\,\boldsymbol{x}) \, dP(\boldsymbol{x}) \, .$$

• The Bayes classifier for the 0/1 loss has the form of a joint mode:

$$\boldsymbol{h}^*(\boldsymbol{x}) = \operatorname*{arg\,max}_{\boldsymbol{y}} P(\boldsymbol{y} \,|\, \boldsymbol{x})$$

- We note that h^{*} is in general not unique. However, the risk of h^{*}, denoted L^{*}(P), is unique, and is called the Bayes risk.
- The regret of h on P(X, Y) is defined as:

$$\operatorname{reg}_{0/1}(\boldsymbol{h}, P) = L_{0/1}(\boldsymbol{h}, P) - L_{0/1}^*(P)$$

• The goal is to train a classifier *h* with a small regret, ideally equal to zero.

We say that the reduction algorithm is consistent if zero regret solution on the reduced problem implies zero regret on the original problem.

We say that the reduction algorithm is consistent if zero regret solution on the reduced problem implies zero regret on the original problem.

• The Label Tree Classifiers ensure consistency when the perfect (zero regret) classifier in each node of the tree guarantees the perfect global classification of a test example.

We say that the reduction algorithm is consistent if zero regret solution on the reduced problem implies zero regret on the original problem.

- The Label Tree Classifiers ensure consistency when the perfect (zero regret) classifier in each node of the tree guarantees the perfect global classification of a test example.
- The conventional training and classification procedures for Label Tree Classifiers do not ensure consistency!

We say that the reduction algorithm is consistent if zero regret solution on the reduced problem implies zero regret on the original problem.

- The Label Tree Classifiers ensure consistency when the perfect (zero regret) classifier in each node of the tree guarantees the perfect global classification of a test example.
- The conventional training and classification procedures for Label Tree Classifiers do not ensure consistency!



Can Label Tree Classifiers be consistent?

Can Label Tree Classifiers be consistent?

Two approaches...

Filter Trees (FT)¹

- Bottom-up learning algorithm to train binary classifiers of the tree.
- Single elimination tournament on the set of label combinations.
- The node classifiers h_{y^i} are trained as before to predict y_{i+1} , where $y^i = (y_1, \dots, y_i)$.
- FT implicitly transforms the underlying distribution P over multi-label examples into a specific distribution $P_{y^i}^{FT}$ over binary examples at each node y^i .
- This transformation for the 0/1 loss filters out all examples that are misclassified by the lower-level classifiers.
- $h_{y^i}(x)$ predicts y_{i+1} given that all classifiers below predict the subsequent labels correctly:

$$h_{y^i}: x \mapsto (y_{i+1} | y_{j+1} = h_{y^j}(x) : j = i+1, \dots, m)$$

• The inference procedure of FT is straight-forward and uses the greedy search (which is sufficient for obtaining consistent predictions).

¹A. Beygelzimer, J. Langford, and P.D. Ravikumar, *Error-correcting tournaments*. In ALT, pp. 247262, 2009









Theorem ²

For all binary classifiers h and distributions P,

$$\operatorname{reg}_{0/1}(\boldsymbol{h}, P) \le m \operatorname{reg}_{0/1}(\boldsymbol{h}, P^{FT})$$

²A. Beygelzimer, J. Langford, and P.D. Ravikumar, *Error-correcting tournaments*. In ALT, pp. 247262, 2009

- Probabilistic binary classifiers on the tree (e.g., logistic regression).
- Each node classifier $h_{\boldsymbol{y}^i}$ delivers estimates of conditional probability:

$$h_{\boldsymbol{y}^i}: \boldsymbol{x} \mapsto P(y_{i+1} = 1 | \boldsymbol{x}, \boldsymbol{y}^i)$$

We denote these estimates by $Q(y_{i+1} = 1 | \boldsymbol{x}, \boldsymbol{y}^i)$.

- Training phase as in the conventional Label Tree Classifiers (can be easily parallelized).
- More complex inference procedure (the greedy approach does not guarantee to find the optimal solution).

Probabilistic Classifier Trees – ϵ -approximate inference

- 1: Input: x, priority lists: $\mathcal{Q} \leftarrow \{y_0\}$, $\mathcal{K} \leftarrow \{\}$, $\epsilon \leftarrow 2^{-c}$ with $c \leq m$ 2: while $\mathcal{Q} \neq \emptyset$ do
- 3: $v \leftarrow \mathsf{pop} \mathsf{ first element in } \mathcal{Q}$
- 4: if v is a leaf then delete all elements in \mathcal{K} and break
- 5: $v_1 \leftarrow (v, 1)$ (left child of v) and $v_0 \leftarrow (v, 0)$ (right child of v)
- 6: compute $Q(\boldsymbol{v}_1 \,|\, \boldsymbol{x})$ and $Q(\boldsymbol{v}_0 \,|\, \boldsymbol{x})$ recursively from $Q(\boldsymbol{v} \,|\, \boldsymbol{x})$
- 7: **if** $Q(\boldsymbol{v}_1 \,|\, \boldsymbol{x}) \geq \epsilon$ **then** add \boldsymbol{v}_1 to \mathcal{Q} sorted in desc. order of Q
- 8: **if** $Q(\boldsymbol{v}_0 \mid \boldsymbol{x}) \geq \epsilon$ **then** add \boldsymbol{v}_0 to \mathcal{Q} sorted in desc. order of Q
- 9: if v_1 and v_0 are not in Q then add v to K in desc. order of Q10: end while
- 11: $\epsilon \leftarrow 0$
- 12: while $\mathcal{K} \neq \emptyset$ do
- 13: $v' \leftarrow \mathsf{pop} \ \mathsf{first} \ \mathsf{element} \ \mathsf{in} \ \mathcal{K}$
- 14: $v' \leftarrow$ apply greedy search downward on v'
- 15: if $Q(v' | x) \ge \epsilon$ then $v \leftarrow v'$ and $\epsilon \leftarrow Q(v' | x)$

16: end while

17: return $v = (y_1, ..., y_m)$

Theorem

For all binary probabilistic classifiers h and distributions the total regret of P, PCT with ϵ -approximate inference:

$$\operatorname{reg}_{0/1}(\boldsymbol{h}, P) \le \sqrt{2m\overline{\operatorname{reg}}_{\log}(h, P)} + 2^{-c} - 2^{-m}$$

Experimental Results

		0/1[%]	t_{train}	t_{test}	А	min	max	1q	2q	3q	P@5
Yeast $m = 14$	BR	85.06	18	0.18	14	14	14	14	14	14	_
	FT	78.30	16	0.17	14	14	14	14	14	14	-
	$PCT_{\epsilon=0.5}$	79.17	21	0.19	14	14	14	14	14	14	-
	$PCT_{\epsilon=0.25}$	77.54	21	0.32	23	14	42	15	26	28	_
	$PCT_{\epsilon=0.0}$	76.34	21	0.37	25	14	61	17	23	30	-
Enron $m = 53$	BR	86.36	66	1.62	53	53	53	53	53	53	-
	FT	82.56	76	1.57	53	53	53	53	53	53	-
	$PCT_{\epsilon=0.5}$	83.77	92	1.49	53	53	53	53	53	53	-
	$PCT_{\epsilon=0.25}$	82.73	92	1.82	72	54	147	54	54	58	_
	$PCT_{\epsilon=0.0}$	81.87	92	3.12	144	53	888	73	106	175	-
$\begin{array}{l} Mediamill \\ m = 101 \end{array}$	BR	92.96	4947	26.20	101	101	101	101	101	101	-
	FT	89.86	7551	27.76	101	101	101	101	101	101	-
	$PCT_{\epsilon=0.5}$	90.33	8695	30.73	101	101	101	101	101	101	-
	$PCT_{\epsilon=0.25}$	90.12	8695	40.86	125	101	260	102	102	102	-
	$PCT_{\epsilon=0.0}$	90.07	8695	54.51	208	101	1756	149	182	228	-
k = 1000	1vsA	91.80	13435	581.43	1000	1000	1000	1000	1000	1000	19.75
	FT	95.77	1051	24.50	10	10	10	10	10	10	-
	$PCT_{\epsilon=0.5}$	96.92	1148	27.18	10	10	10	10	10	10	3.08
	$PCT_{\epsilon=0.25}$	94.99	1148	43.60	22	10	30	21	21	21	5.55
	$PCT_{\epsilon=0.0}$	92.93	1148	69.56	55	10	192	38	52	69	15.16
	$PCT'_{\epsilon=0.0}$	92.66	2379	109.55	94	17	237	75	91	111	19.75

Conclusions

Filter Trees

- Logarithmic prediction time in the number of classes or label combinations.
- FT can be used with any type of binary classifiers.
- To guarantee the consistency of the greedy prediction it requires more demanding training.
- Filtering may reduce the number of training examples in the top levels of the tree.
- There is no option to predict top classes with the highest conditional probabilities.

Probabilistic Classifier Trees

- The prediction can be longer than for FT.
- The complexity of the classification is connected with the noise of data.
- PCT requires the use of probabilistic classifiers.
- The learning is much simpler and can be easily parallelized.
- The probabilistic nature allows to deliver a list of top classes.

Thank you!

Questions?

For more check: www.cs.put.poznan.pl/ajachnik or contact me: ajachnik@cs.put.poznan.pl

This project is partially supported by the Foundation of Polish Science under the Homing Plus programme, co-financed by the European Regional Development Fund.





