

Transductive Bounds for the Multi-class Majority Vote Classifier

Vasilii Feofanov, Emilie Devijver, Massih-Reza Amini

University Grenoble Alpes, Grenoble INP
LIG, CNRS, Grenoble 38000, France
(firstname.lastname@univ-grenoble-alpes.fr)

AAAI, 2019

Introduction

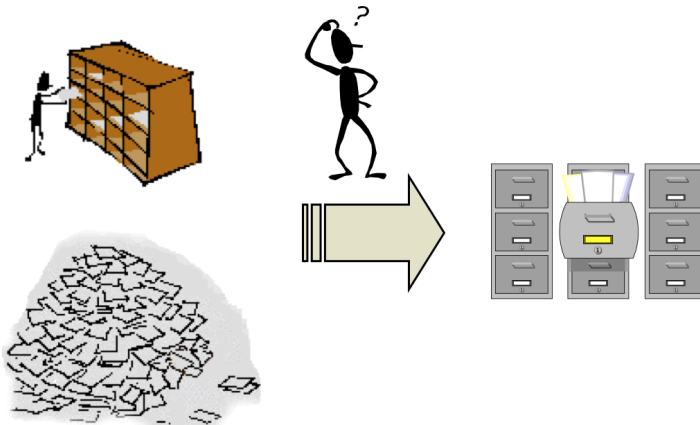
Framework

Transductive
Bounds

Application

Introduction

In many applications, labeling examples is prohibitive while huge number of unlabeled data are available.



Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

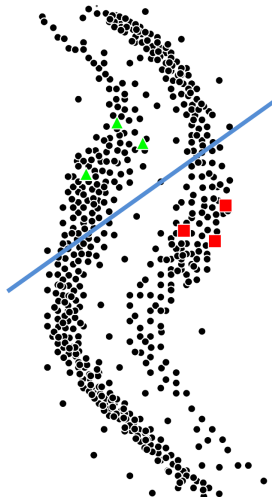
Framework

Transductive
Bounds

Application

Introduction

In many applications, labeling examples is prohibitive while huge number of unlabeled data are available.



Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

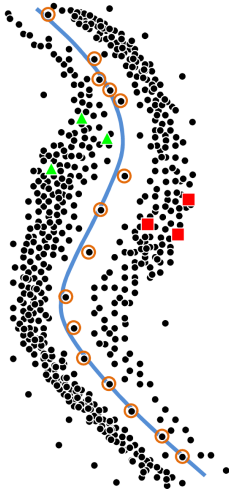
Framework

Transductive
Bounds

Application

Introduction

In many applications, labeling examples is prohibitive while huge number of unlabeled data are available. Taking into account the margin to find the low density regions of labeled and unlabeled examples constitutes the basis of many semi-supervised learning algorithms.



Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application

- To increase performance on the unlabelled set we can infer a model in the transductive way.
- The PAC-Bayesian theory [McAllester, 1999] proposes transductive risk bounds for the Gibbs and the Bayes classifiers.
- Based on this, the self-learning algorithm has been proposed for the binary classification [Amini et al., 2008]. It iteratively pseudo-labels unlabelled examples depending on the prediction scores.

1. Multi-class PAC-Bayesian Theory concerns only the **inductive** case by now.
 - No **transductive** bound of the Bayes classifier has been proposed yet.
2. Few **multi-class** SSL approaches.
3. No self-learning algorithm for the multi-class case yet.

Problem statement

- Input space $\mathcal{X} \subset \mathbb{R}^d$, output $\mathcal{Y} = \{1, \dots, K\}$ space.
- Labelled examples $Z_{\mathcal{L}} = \{(\mathbf{x}_i, y_i)\}_{i=1}^l$.
- Unlabelled observations $X_{\mathcal{U}} = \{\mathbf{x}_i\}_{i=l+1}^{l+u}$.
- $(\mathbf{x}_i, y_i) \in Z_{\mathcal{L}} \sim \mathcal{D}$ i.i.d.
- Marginal distribution P_X defined over \mathcal{X} .
- Assumption: $\forall \mathbf{x}_i \in X_{\mathcal{U}}$, there is exactly one possible label.
- Hypothesis space \mathcal{H} .
- Prior P and posterior Q distributions over \mathcal{H} .

Goal: accurate classification of the unlabelled set.

Context: $l \ll u$.

Definitions

- $B_Q(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} [\mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}]$, – Bayes majority
vote classifier
- $G_Q(\mathbf{x}) := \operatorname{rand}_{h \sim Q} h(\mathbf{x})$, – Gibbs stochastic classifier

Definitions

- $B_Q(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} [\mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}]$,
- $G_Q(\mathbf{x}) := \operatorname{rand}_{h \sim Q} h(\mathbf{x})$,
- $R_{\mathcal{U}}(B_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{B_Q(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $R_{\mathcal{U}}(G_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,

The error to **predict j** given **class i**.

- $B_Q(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} [\mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}]$,
- $G_Q(\mathbf{x}) := \operatorname{rand}_{h \sim Q} h(\mathbf{x})$,
- $R_{\mathcal{U}}(B_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{B_Q(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $R_{\mathcal{U}}(G_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $E_{\mathcal{U}}(h) := \frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{h(\mathbf{x}') \neq y'}$, – error rate
- $\mathbf{C}_h^{\mathcal{U}} := (R_{\mathcal{U}}(h, i, j))_{i,j=\{1,\dots,K\}^2, i \neq j}$, – confusion matrix¹

¹[Morvant et al., 2012]

- $B_Q(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} [\mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}]$,
- $G_Q(\mathbf{x}) := \operatorname{rand}_{h \sim Q} h(\mathbf{x})$,
- $R_{\mathcal{U}}(B_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{B_Q(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $R_{\mathcal{U}}(G_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $E_{\mathcal{U}}(h) := \frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{h(\mathbf{x}') \neq y'}$,
- $\mathbf{C}_h^{\mathcal{U}} := (R_{\mathcal{U}}(h, i, j))_{\substack{i,j=\{1,\dots,K\} \\ i \neq j}}$,
- $m_Q(\mathbf{x}, c) = \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}$, – margin: indicator of confidence

- $B_Q(\mathbf{x}) := \operatorname{argmax}_{c \in \mathcal{Y}} [\mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}]$,
- $G_Q(\mathbf{x}) := \operatorname{rand}_{h \sim Q} h(\mathbf{x})$,
- $R_{\mathcal{U}}(B_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{B_Q(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $R_{\mathcal{U}}(G_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x}')=j} \mathbb{1}_{y'=i}$,
- $E_{\mathcal{U}}(h) := \frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{h(\mathbf{x}') \neq y'}$,
- $\mathbf{C}_h^{\mathcal{U}} := (R_{\mathcal{U}}(h, i, j))_{\substack{i,j=\{1,\dots,K\} \\ i \neq j}}^2$,
- $m_Q(\mathbf{x}, c) = \mathbb{E}_{h \sim Q} \mathbb{1}_{h(\mathbf{x})=c}$,
- $R_{\mathcal{U} \wedge \theta}(B_Q, i, j) := \frac{1}{u_i} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{B_Q(\mathbf{x}')=j} \mathbb{1}_{y'=i} \mathbb{1}_{m_Q(\mathbf{x}', j) \geq \theta_j}$,
 - risk to have the conditional error and the margin above θ_j

A Transductive Bound for the Conditional Risk

Theorem

$\forall Q$ and $\forall \delta \in (0, 1]$, $\forall \theta \in [0, 1]^K$ with prob. at least $1 - \delta$:

$$R_u(B_Q, i, j) \leq \inf_{\gamma \in [0, 1]} \left\{ I_{i,j}^{(\leq, <)}(0, \gamma) + \frac{1}{\gamma} \left[(K_{i,j}^\delta - M_{i,j}^{<}(\gamma)) \right]_+ \right\},$$

$$R_{u \wedge \theta}(B_Q, i, j) \leq \inf_{\gamma \in [\theta_j, 1]} \left\{ I_{i,j}^{(\leq, <)}(\theta_j, \gamma) + \frac{1}{\gamma} \left[(K_{i,j}^\delta - M_{i,j}^{<}(\gamma) + M_{i,j}^{<}(\theta_j)) \right]_+ \right\},$$

where

- $K_{i,j}^\delta = R_u^\delta(G_Q, i, j) - \varepsilon_{i,j}$,
- $R_u^\delta(G_Q, i, j)$ is an upper bound that holds with prob. at least $1 - \delta$.
- $\varepsilon_{i,j}$ is the average of j -margins in class i and class j is not predicted,
- $I_{i,j}^{(<1, <2)}(t, s)$ is the proportion of class i examples with the margin between t and s ,
- $M_{i,j}^{<}(t)$ is the average of j -margins in class i that less than t .

Proof of Theorem

- Lemma 1: connection between the Gibbs conditional risk and the joint Bayes one.
- Bound derived from a solution of a linear program where the error is maximized.
- Lemma 2: the solution of the linear program is the maximal feasible solution in the lexicographical order.

Proposition

Ass.: $\forall \mathbf{x}' \in X_{\mathcal{U}} \exists C \in [0, 1]$ s.t. $\forall (i, j) \in \mathcal{Y}^2, \forall \gamma > 0$ if B_Q makes cond. mistakes on examples with the margin $\gamma \Rightarrow$ the proportion of cond. misclassified examples with margin $< \gamma$ is lower bounded by C . Then, with prob. at least $1 - \delta$:

$$F_{i,j}^{\delta} - R_{\mathcal{U}}(B_Q, i, j) \leq \frac{1 - C}{C} R_{\mathcal{U}}(B_Q, i, j) + \frac{R_u^{\delta}(G_Q, i, j) - R_{\mathcal{U}}(G_Q, i, j)}{\gamma^*},$$

where $F_{i,j}^{\delta}$ is the proposed bound, and γ^* is max margin on which the Bayes classifier makes a cond. mistake.

Proposition

Ass.: $\forall \mathbf{x}' \in X_{\mathcal{U}} \exists C \in [0, 1]$ s.t. $\forall (i, j) \in \mathcal{Y}^2, \forall \gamma > 0$ if B_Q makes cond. mistakes on examples with the margin $\gamma \Rightarrow$ the proportion of cond. misclassified examples with margin $< \gamma$ is lower bounded by C . Then, with prob. at least $1 - \delta$:

$$F_{i,j}^{\delta} - R_{\mathcal{U}}(B_Q, i, j) \leq \frac{1 - C}{C} R_{\mathcal{U}}(B_Q, i, j) + \frac{R_u^{\delta}(G_Q, i, j) - R_{\mathcal{U}}(G_Q, i, j)}{\gamma^*},$$

where $F_{i,j}^{\delta}$ is the proposed bound, and γ^* is max margin on which the Bayes classifier makes a cond. mistake.

Additional assumptions:

- The Gibbs conditional risk bound is tight,

Proposition

Ass.: $\forall \mathbf{x}' \in X_{\mathcal{U}} \exists C \in [0, 1]$ s.t. $\forall (i, j) \in \mathcal{Y}^2, \forall \gamma > 0$ if B_Q makes cond. mistakes on examples with the margin $\gamma \Rightarrow$ the proportion of cond. misclassified examples with margin $< \gamma$ is lower bounded by C . Then, with prob. at least $1 - \delta$:

$$F_{i,j}^{\delta} - R_{\mathcal{U}}(B_Q, i, j) \leq \frac{1-C}{C} R_{\mathcal{U}}(B_Q, i, j) + \frac{R_u^{\delta}(G_Q, i, j) - R_{\mathcal{U}}(G_Q, i, j)}{\gamma^*},$$

where $F_{i,j}^{\delta}$ is the proposed bound, and γ^* is max margin on which the Bayes classifier makes a cond. mistake.

Additional assumptions:

- The Gibbs conditional risk bound is tight,
- The Bayes classifier makes its mistakes mostly on examples with low margins $\Rightarrow C$ is close to 1.

Proposition

Ass.: $\forall \mathbf{x}' \in X_{\mathcal{U}} \exists C \in [0, 1]$ s.t. $\forall (i, j) \in \mathcal{Y}^2, \forall \gamma > 0$ if B_Q makes cond. mistakes on examples with the margin $\gamma \Rightarrow$ the proportion of cond. misclassified examples with margin $< \gamma$ is lower bounded by C . Then, with prob. at least $1 - \delta$:

$$F_{i,j}^{\delta} - R_{\mathcal{U}}(B_Q, i, j) \leq \frac{1-C}{C} R_{\mathcal{U}}(B_Q, i, j) + \frac{R_{\mathcal{U}}^{\delta}(G_Q, i, j) - R_{\mathcal{U}}(G_Q, i, j)}{\gamma^*},$$

where $F_{i,j}^{\delta}$ is the proposed bound, and γ^* is max margin on which the Bayes classifier makes a cond. mistake.

Additional assumptions:

- The Gibbs conditional risk bound is tight,
- The Bayes classifier makes its mistakes mostly on examples with low margins $\Rightarrow C$ is close to 1.

Hence, the bound is tight!

A Transductive Bound of the Error Rate

Corollary

$U_{i,j}^{\delta}(\theta)$ is a bound of $R_{\mathcal{U} \wedge \theta}(B_Q, i, j)$ holding with prob. at least $1 - \delta$,

$\mathbf{U}_{\theta}^{\delta}$ the corresponding confusion matrix.

Then, we have:

$$\begin{aligned} \mathbb{E}_{\mathcal{U} \wedge \theta}(B_Q) &\leq \left\| \left(\mathbf{U}_{\theta}^{\delta} \right)^{\top} \mathbf{p} \right\|_1, \\ \mathbb{E}_{\mathcal{U}}(B_Q) &\leq \left\| \left(\mathbf{U}_{\mathbf{0}_K}^{\delta} \right)^{\top} \mathbf{p} \right\|_1, \end{aligned}$$

where $\mathbf{p} = \{u_i/u\}_{i=1}^K$.

Automatic Threshold Finding

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

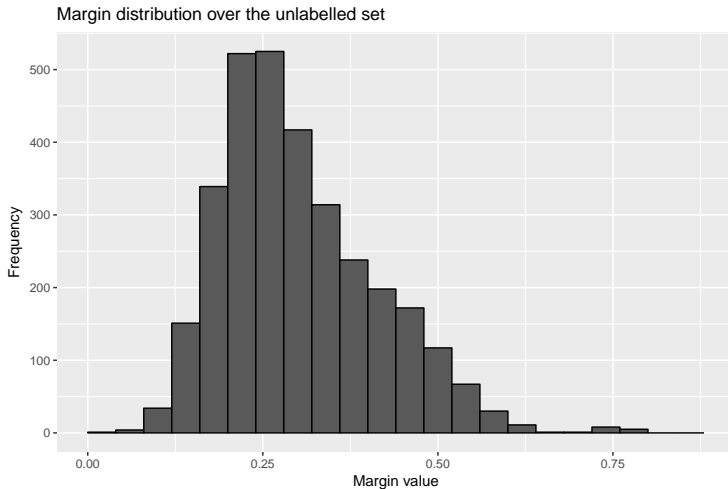
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Automatic Threshold Finding

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

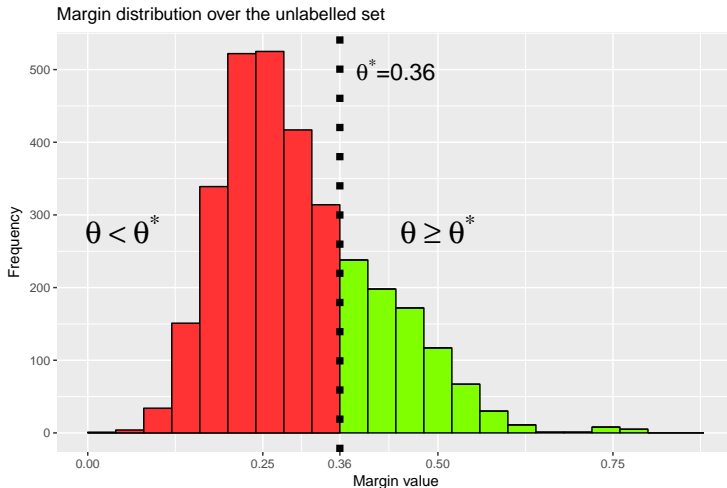
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Definition

Conditional Bayes error rate $E_{\mathcal{U}|\theta}(B_Q)$:

$$E_{\mathcal{U}|\theta}(B_Q) := \frac{E_{\mathcal{U} \wedge \theta}(B_Q)}{\frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{m_Q(\mathbf{x}', B_Q(\mathbf{x}')) \geq \theta_{B_Q}(\mathbf{x}')}},$$

Trade-off between:

Introduction

Framework

Transductive
Bounds

Application

Definition

Conditional Bayes error rate $E_{\mathcal{U}|\theta}(B_Q)$:

$$\mathbb{E}_{\mathcal{U}|\theta}(B_Q) := \frac{\mathbb{E}_{\mathcal{U} \wedge \theta}(B_Q)}{\frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{m_Q(\mathbf{x}', B_Q(\mathbf{x}')) \geq \theta_{B_Q(\mathbf{x}')}}},$$

Trade-off between:

- Error on unlabeled examples with margin above $\theta_{B_Q(\mathbf{x}')}$,

Conditional Bayes Error

Definition

Conditional Bayes error rate $E_{\mathcal{U}|\theta}(B_Q)$:

$$E_{\mathcal{U}|\theta}(B_Q) := \frac{E_{\mathcal{U} \wedge \theta}(B_Q)}{\frac{1}{u} \sum_{\mathbf{x}' \in X_{\mathcal{U}}} \mathbb{1}_{m_Q(\mathbf{x}', B_Q(\mathbf{x}')) \geq \theta_{B_Q}(\mathbf{x}')}},$$

Trade-off between:

- Error on unlabeled examples with margin above $\theta_{B_Q}(\mathbf{x}')$,
- Fraction of pseudo-labeled examples in $X_{\mathcal{U}}$.

Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application

Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

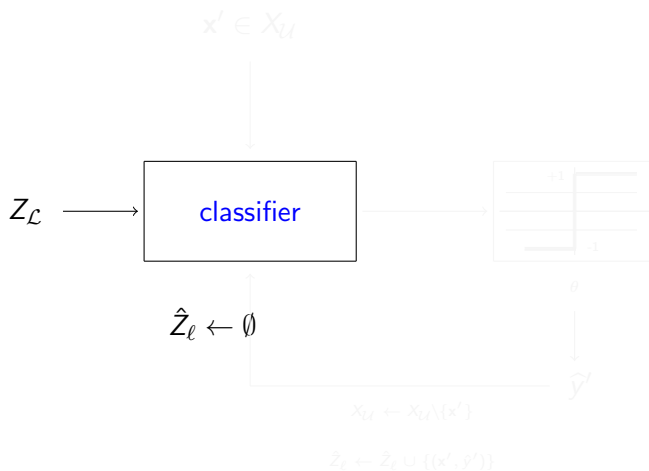
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

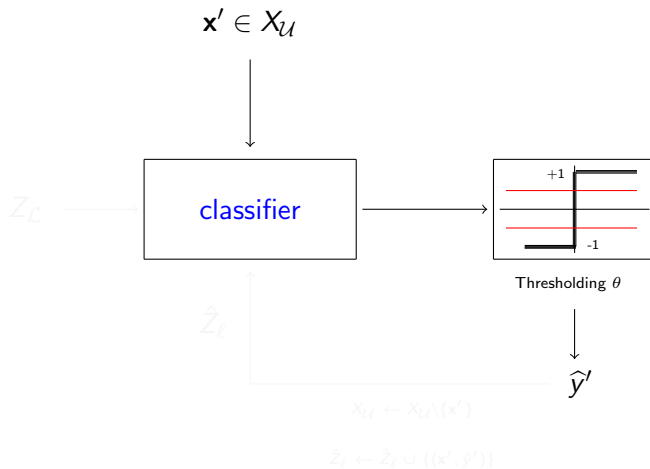
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

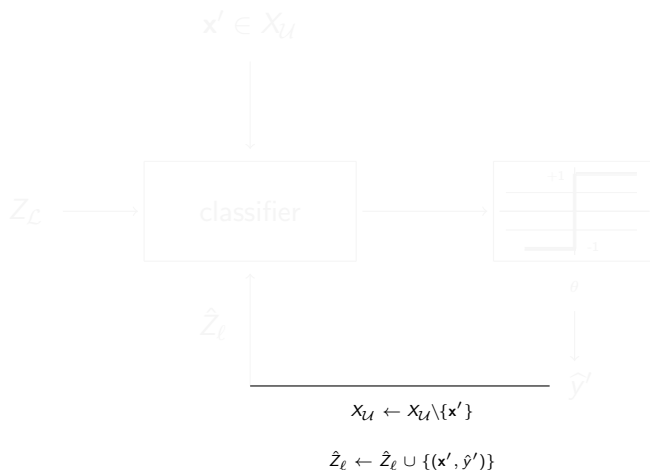
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

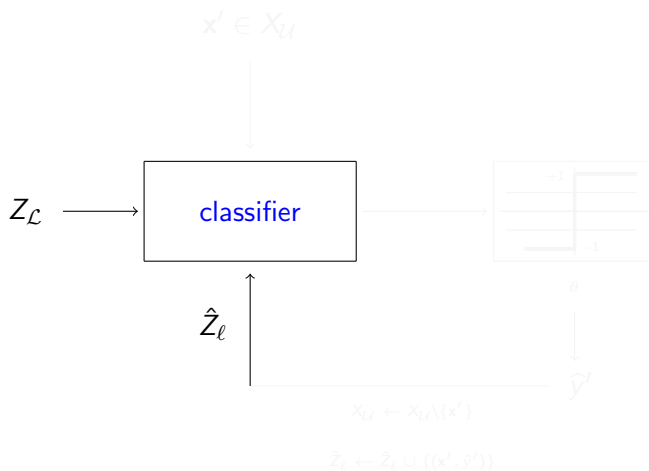
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Multi-class Self-learning Algorithm

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

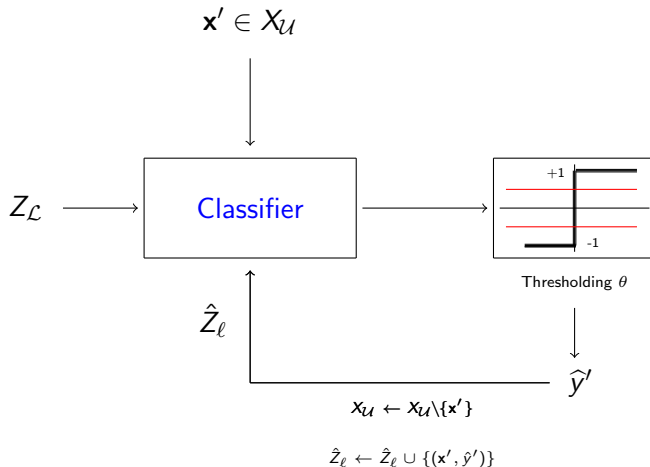
Vasilii Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application



Classification Performance with respect to l .

Transductive
Bounds for the
Multi-class
Majority Vote
Classifier

Vasili Feofanov,
Emilie Devijver,
Massih-Reza Amini

Introduction

Framework

Transductive
Bounds

Application

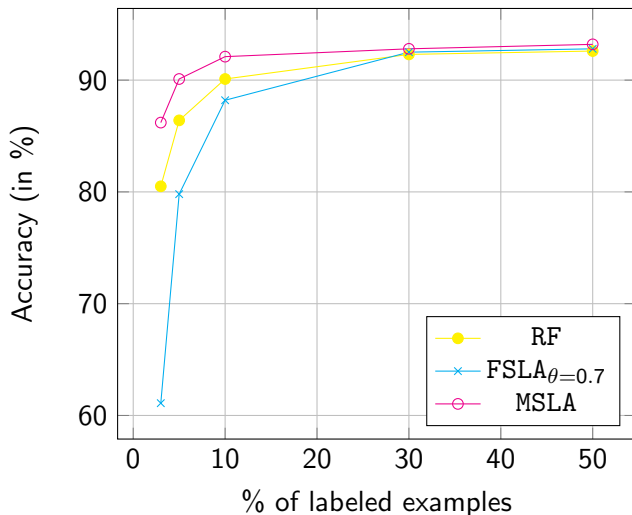


Figure: Accuracy on 3500 examples randomly chosen from the MNIST dataset.

Conclusion and Perspectives

- Proposed transductive bounds for the Bayes classifier, which are tight under certain conditions.
- Self-learning with automatic threshold finding shows promising results for semi-supervised tasks.

Conclusion and Perspectives

- Proposed transductive bounds for the Bayes classifier, which are tight under certain conditions.
- Self-learning with automatic threshold finding shows promising results for semi-supervised tasks.
- Future perspective: self-learning with semi-supervised feature selection.

