

Lecture 8: Discriminative structured output learning, Perceptron algorithm

Vojtěch Franc

April 16, 2015

- 7.A: Definition of structured classification task and its solution via generative and discriminative learning
- 7.B: Implementation of ERM learning using Perceptron algorithm
- 7.C: Learning of max-sum classifier

XEP33SML – Structured Model Learning, Summer 2015

7.A: Structured Output Classification

The setting

- ◆ $\mathbf{x} \in \mathcal{X}$ is input (also observation or measurement)
- ◆ $\mathbf{y} \in \mathcal{Y}$ is output (also hidden state or label)
- ◆ (\mathbf{x}, \mathbf{y}) are assumed to be realizations of random variables (X, Y) with p.d.f. $p(\mathbf{x}, \mathbf{y})$
- ◆ $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow [0, \infty)$ loss function

The task is to find the optimal (Bayes) classification strategy $h: \mathcal{X} \rightarrow \mathcal{Y}$ achieving the minimal expected risk

$$R^* = \inf_{h: \mathcal{X} \rightarrow \mathcal{Y}} R(h) \quad \text{where} \quad R(h) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p} [\ell(\mathbf{y}, h(\mathbf{x}))]$$

Taxonomy based on the set \mathcal{Y} :

1. (Flat) classification: $\mathcal{Y} = \{1, \dots, K\}$

2. Structured output classification: \mathcal{Y} is a finite set of structured objects

E.g. $\mathbf{y} = (y_i \in \mathcal{Y}_i \mid i \in \mathcal{V}) \in \mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_{|\mathcal{V}|}$ is a labeling of a set of parts \mathcal{V} .

7.A: Learning: using examples to solve the task

- ◆ Typically, $p(\mathbf{x}, \mathbf{y})$ is unknown but we are given training set

$$\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \mathcal{I} = \{1, \dots, m\}\}$$

drawn from i.i.d. random variables with p.d.f. $p(\mathbf{x}, \mathbf{y})$.

- ◆ Given training set \mathcal{T} , the learning algorithm $A: \cup_{m=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathcal{H}$ selects a classifier $h_{\mathcal{T}}$ from a fixed class of strategies $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$.
- ◆ The learning algorithm is characterized by the excess error:

$$\underbrace{\left(R(h_{\mathcal{T}}) - R^* \right)}_{\text{excess error}} = \underbrace{\left(R(h_{\mathcal{T}}) - R(h_{\mathcal{H}}) \right)}_{\text{estimation error}} + \underbrace{\left(R(h_{\mathcal{H}}) - R^* \right)}_{\text{approximation error}}$$

where $h_{\mathcal{H}}$ denotes the optimal classifier within the class \mathcal{H} , i.e

$$R(h_{\mathcal{H}}) = \inf_{h \in \mathcal{H}} R(h)$$

- ◆ Selection of the class of strategies \mathcal{H} controls the trade-off between the estimation error (depends on \mathcal{H} and m) and the approximation error (depends on \mathcal{H}).

7.A: Generative learning (lectures 3-7)

1. Use the training set $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \{1, \dots, m\}\}$ to approximate $p(\mathbf{x}, \mathbf{y})$ by $\hat{p}(\mathbf{x}, \mathbf{y})$.

For example, use the Maximum-Likelihood method:

- (a) Guess the shape of the distribution, e.g. $\hat{p}_{\mathbf{u}}(\mathbf{x}, \mathbf{y}) = \frac{1}{Z(\mathbf{u})} \exp\langle \mathbf{u}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$, $\mathbf{u} \in \mathcal{U}$
- (b) Find the ML estimate

$$\mathbf{u}_{\mathcal{T}} \in \operatorname{argmax}_{\mathbf{u} \in \mathcal{U}} \sum_{i=1}^m \log \hat{p}_{\mathbf{u}}(\mathbf{x}^i, \mathbf{y}^i)$$

2. Construct a plug-in classifier

$$h_{\mathcal{T}}(\mathbf{x}) \in \operatorname{argmin}_{h: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \hat{p}}[\ell(\mathbf{y}, h(\mathbf{x}))]$$

7.A: Discriminative learning

- ◆ Replace $R(h)$ by a surrogate risk $R_{\mathcal{T}}(h)$ which can be evaluated on the training set $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \{1, \dots, m\}\}$ and find the classifier by solving

$$h_{\mathcal{T}} \in \operatorname{argmin}_{h \in \mathcal{H}} R_{\mathcal{T}}(h)$$

where $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ is a fixed class of strategies.

- ◆ For example, the **Empirical Risk Minimization** (ERM) learning uses the **empirical risk**

$$R_{\mathcal{T}}(h) = \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{y}^i, h(\mathbf{x}^i))$$

7.A: Statistical Consistency

Definition 1. Let $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ be a class of strategies. Let $A: \cup_{m=1}^{\infty} (\mathcal{X} \times \mathcal{Y})^m \rightarrow \mathcal{H}$ be a learning algorithm which returns a classifier $h_{\mathcal{T}} \in \mathcal{H}$ for given training set $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \{1, \dots, m\}\}$ generated from i.i.d. with $p(\mathbf{x}, \mathbf{y})$. The algorithm A is **statistically consistent** in \mathcal{H} w.r.t. $p(\mathbf{x}, \mathbf{y})$ if for all $\varepsilon > 0$ it holds that

$$\lim_{m \rightarrow \infty} \mathbb{P}_{\mathcal{T} \sim p^m} \left(R(h_{\mathcal{T}}) - R(h_{\mathcal{H}}) \geq \varepsilon \right) = 0$$

If A is consistent for all $p(\mathbf{x}, \mathbf{y})$ then A is so called **universally consistent** in \mathcal{H} .

In words, with the number of examples going to infinity the estimation error **converges in probability to** zero.

7.A: Consistency of ERM

The estimation error of ERM can be upper bounded as follows:

$$\begin{aligned} R(h_{\mathcal{T}}) - R(h_{\mathcal{H}}) &= \left(R(h_{\mathcal{T}}) - R_{\mathcal{T}}(h_{\mathcal{T}}) \right) + \underbrace{\left(R_{\mathcal{T}}(h_{\mathcal{T}}) - R_{\mathcal{T}}(h_{\mathcal{H}}) \right)}_{\text{non-pos. due to ERM}} + \left(R_{\mathcal{T}}(h_{\mathcal{H}}) - R(h_{\mathcal{H}}) \right) \\ &\leq \left(R(h_{\mathcal{T}}) - R_{\mathcal{T}}(h_{\mathcal{T}}) \right) + \left(R_{\mathcal{T}}(h_{\mathcal{H}}) - R(h_{\mathcal{H}}) \right) \\ &\leq 2 \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}}(h) \right| \end{aligned}$$

Therefore $\varepsilon \leq R(h_{\mathcal{T}}) - R(h_{\mathcal{H}})$ implies $\frac{\varepsilon}{2} \leq \sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}}(h) \right|$ and thus

$$\mathbb{P}_{\mathcal{T} \sim p^m} \left(R(h_{\mathcal{T}}) - R(h_{\mathcal{H}}) \geq \varepsilon \right) \leq \mathbb{P}_{\mathcal{T} \sim p^m} \left(\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}}(h) \right| \geq \frac{\varepsilon}{2} \right)$$

7.A: Uniform convergence of empirical risk

Definition 2. *The empirical risk converges uniformly to the expected risk in $\mathcal{H} \subseteq \mathcal{Y}^{\mathcal{X}}$ if for all $\varepsilon > 0$ it holds that*

$$\lim_{m \rightarrow \infty} \mathbb{P}_{\mathcal{T} \sim p^m} \left(\sup_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}}(h) \right| \geq \varepsilon \right) = 0$$

Corrolary: The uniform convergence of the empirical risk in \mathcal{H} implies the consistency of ERM in \mathcal{H} as we see from previous slide.

7.A: Uniform convergence of empirical risk

Two examples of \mathcal{H} when the empirical risk converges uniformly:

1. \mathcal{H} is a finite set. Then,

$$\mathbb{P}_{\mathcal{T} \sim p^m} \left(\max_{h \in \mathcal{H}} \left| R(h) - R_{\mathcal{T}}(h) \right| \geq \varepsilon \right) \leq 2|\mathcal{H}| \exp \left(\frac{-2m \varepsilon^2}{l_{max}} \right)$$

where $l_{max} = \max_{\mathbf{y}, \mathbf{y}'} \ell(\mathbf{y}, \mathbf{y}')$.

A similar bound for the ML learning was explained in Lecture 6, section 5.

2. $\ell(\mathbf{y}, \mathbf{y}') = [\mathbf{y} \neq \mathbf{y}']$ and VC-dimension of \mathcal{H} is finite.

VC-dimension (capacity) of \mathcal{H} is the maximal number of inputs which can be classified by strategies from \mathcal{H} in all possible (that is $|\mathcal{Y}|^d$) ways.

7.B: Learning linear classifier from separable examples

Let us consider \mathcal{H} composed of all linear classifiers

$$h(\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$$

where $\mathbf{w} \in \mathbb{R}^n$ are the parameters and $\Psi: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^n$ is a fixed function (feature map).

Definition 3. (*Linearly separable examples*) The examples $\{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \mathcal{I} = \{1, \dots, m\}\}$ are called *linearly separable* if there exists $\mathbf{w} \in \mathbb{R}^n$ such that

$$\langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) \rangle > \langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}) \rangle, \quad \forall i \in \mathcal{I}, \mathbf{y} \in \mathcal{Y} \setminus \{\mathbf{y}^i\}$$

Implementation of ERM: For a loss function $\ell(\mathbf{y}, \mathbf{y}) = 0, \mathbf{y} \in \mathcal{Y}$, the linearly separable examples imply that there exists a linear classifier given by $\mathbf{w} \in \mathbb{R}^n$ with zero empirical risk, i.e.,

$$h(\mathbf{x}^i; \mathbf{w}) = \mathbf{y}^i, \quad \forall i \in \mathcal{I}$$

hence the ERM can be translated to solving a set of $m(|\mathcal{Y}| - 1)$ linear inequalities.

7.B: Perceptron algorithm

Task: given a set of point $\{\mathbf{a}^i \in \mathbb{R}^n \mid i \in \mathcal{I}\}$ we want to find $\mathbf{w} \in \mathbb{R}^n$ which satisfies the linear inequalities

$$\langle \mathbf{w}, \mathbf{a}^i \rangle > 0, \quad \forall i \in \mathcal{I}.$$

If the task has a solution then the points are **linearly separable**.

Algorithm 1 Perceptron

1: Set $\mathbf{w}^0 = \mathbf{0}$, $t \leftarrow 0$

2: $t \leftarrow t + 1$

3: Find index $j \in \mathcal{I}$ of a violating inequality such that

$$\langle \mathbf{w}^t, \mathbf{a}^j \rangle \leq 0$$

4: If there is no violating inequality return \mathbf{w}^t . Otherwise update

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \mathbf{a}^j$$

and go to Step 2.

7.B: Convergence of Perceptron

Theorem 1. For any linearly separable points $\{\mathbf{a}^i \in \mathbb{R}^n \mid i \in \mathcal{I}\}$, the Perceptron algorithm terminates in

$$\frac{A^2}{\gamma^2}$$

steps at most where

$$A = \max_{i \in \mathcal{I}} \|\mathbf{a}^i\|_2 \quad \text{and} \quad \gamma = \max_{\mathbf{w} \in \mathbb{R}^n} \min_{i \in \mathcal{I}} \frac{\langle \mathbf{w}, \mathbf{a}^i \rangle}{\|\mathbf{w}\|_2}$$

Note that the upper bound $\frac{A^2}{\gamma^2}$ does not depend on the number of points $m = |\mathcal{I}|$.

7.B: Perceptron applied for learning linear classifier

Task: Given linearly separable training set $\{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \mathcal{I} = \{1, \dots, m\}\}$ find parameters $\mathbf{w} \in \mathbb{R}^n$ of linear classifier

$$h(\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{w}, \Psi(\mathbf{x}, \mathbf{y}) \rangle$$

such that $h(\mathbf{x}^i; \mathbf{w}) = \mathbf{y}^i, i \in \mathcal{I}$.

Algorithm 2 Perceptron

- 1: Set $\mathbf{w} \leftarrow \mathbf{0}$
- 2: Find erroneous example $j \in \mathcal{I}$ such that

$$\mathbf{y}^j \neq \hat{\mathbf{y}}^j \quad \text{where} \quad \hat{\mathbf{y}}^j = h(\mathbf{x}^j)$$

- 3: If there is no erroneous example return \mathbf{w} . Otherwise update

$$\mathbf{w} \leftarrow \mathbf{w} + \Psi(\mathbf{x}^j, \mathbf{y}^j) - \Psi(\mathbf{x}^j, \hat{\mathbf{y}}^j)$$

and go to Step 2.

7.B: Perceptron applied for learning linear classifier

- ◆ By Theorem 1 we have a guarantee that for linearly separable training set $\{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X} \times \mathcal{Y} \mid i \in \mathcal{I}\}$ the Perceptron terminates after at most $\frac{A^2}{\gamma^2}$ iterations where

$$A = \max_{i \in \mathcal{I}, \mathbf{y} \in \mathcal{Y} \setminus \{\mathbf{y}^i\}} \|\Psi(\mathbf{x}^i, \mathbf{y}^i) - \Psi(\mathbf{x}^i, \mathbf{y})\|$$

and

$$\gamma = \max_{\mathbf{w} \in \mathbb{R}^n} \min_{i \in \mathcal{I}, \mathbf{y} \in \mathcal{Y} \setminus \{\mathbf{y}^i\}} \frac{\langle \mathbf{w}, \Psi(\mathbf{x}^i, \mathbf{y}^i) - \Psi(\mathbf{x}^i, \mathbf{y}) \rangle}{\|\mathbf{w}\|_2}$$

- ◆ Perceptron requires only the classification oracle $h(\mathbf{x}; \mathbf{w})$ which does the hard task of searching for the violating inequality among $|\mathcal{Y}|$ options.

7.C: Max-sum classifier

Setting:

- ◆ $(\mathcal{V}, \mathcal{E})$ is undirected graph; \mathcal{V} are parts and $\mathcal{E} \subseteq \binom{|\mathcal{V}|}{2}$ pairs of related parts
- ◆ each part $v \in \mathcal{V}$ described by observation $x \in \mathcal{X}$ and label $y \in \mathcal{Y}$; \mathcal{X} and \mathcal{Y} are finite
- ◆ $q_v: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ quality of label y_v given x_v ; $\mathbf{q} = (q_v(x, y) \in \mathbb{R} \mid x \in \mathcal{X}, y \in \mathcal{Y}, v \in \mathcal{V})$
- ◆ $g_{vv'}: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ quality a label pair $(y_v, y_{v'})$;
 $\mathbf{g} = (g_{vv'}(y, y') \in \mathbb{R} \mid (y, y') \in \mathcal{Y}^2, \{v, v'\} \in \mathcal{E})$

Max-sum classifier: Given observations $\mathbf{x} = (x_v \in \mathcal{X} \mid v \in \mathcal{V}) \in \mathcal{X}^{\mathcal{V}}$, the max-sum classifier $h: \mathcal{X}^{\mathcal{V}} \rightarrow \mathcal{Y}^{\mathcal{V}}$ returns labeling $\mathbf{y} = (y_v \in \mathcal{Y} \mid v \in \mathcal{V}) \in \mathcal{Y}^{\mathcal{V}}$ with the maximal overall quality

$$h(\mathbf{x}; \mathbf{q}, \mathbf{g}) \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} f(\mathbf{x}, \mathbf{y}; \mathbf{q}, \mathbf{g})$$

where

$$f(\mathbf{x}, \mathbf{y}; \mathbf{q}, \mathbf{g}) = \sum_{v \in \mathcal{V}} q_v(x_v, y_v) + \sum_{\{v, v'\} \in \mathcal{E}} g_{vv'}(y_v, y_{v'})$$

The max-sum classifier is an instance of the linear classifier since $f(\mathbf{x}, \mathbf{y}; \mathbf{q}, \mathbf{g}) = \langle \Psi(\mathbf{x}, \mathbf{y}), \mathbf{w} \rangle$ where $\mathbf{w} = (\mathbf{q}, \mathbf{g})$ and $\Psi: \mathcal{X}^{\mathcal{V}} \times \mathcal{Y}^{\mathcal{V}} \rightarrow \mathbb{R}^{|\mathcal{Y}| \cdot |\mathcal{V}| + |\mathcal{E}| \cdot |\mathcal{Y}|^2}$ is constructed appropriately.

7.C: Relation between Max-sum classifier and Gibbs distribution

- ◆ $(\mathcal{V}, \mathcal{E})$ is undirected graph
- ◆ $\{(X_v, Y_v) \mid v \in \mathcal{V}\}$ is a field of random variables taking values from $(x_v, y_v) \in X \times \mathcal{Y}, v \in \mathcal{V}$
- ◆ the random variables are distributed according to the Gibbs distribution

$$\begin{aligned} p_{\mathbf{q}, \mathbf{g}}(\mathbf{x}, \mathbf{y}) &= \frac{1}{Z(\mathbf{q}, \mathbf{g})} \exp \left(\sum_{v \in \mathcal{V}} q_v(x_v, y_v) + \sum_{\{v, v'\} \in \mathcal{E}} g_{vv'}(y_v, y_{v'}) \right) \\ &= \frac{1}{Z(\mathbf{q}, \mathbf{g})} \exp f(\mathbf{x}, \mathbf{y}; \mathbf{q}, \mathbf{g}) \end{aligned}$$

- ◆ The optimal (Bayes) classifier minimizing the expected risk under the 0/1-loss

$$R(h) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\mathbf{q}, \mathbf{g}}} [\mathbf{y} \neq h(\mathbf{x})]$$

is the max-sum classifier

$$h(\mathbf{x}; \mathbf{q}, \mathbf{g}) \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} f(\mathbf{x}, \mathbf{y}; \mathbf{q}, \mathbf{g})$$

7.C: Learning max-sum classifier from linearly separable examples



Task: Given linearly separable training set $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{y}^i) \in \mathcal{X}^{\mathcal{V}} \times \mathcal{Y}^{\mathcal{V}} \mid i \in \mathcal{I} = \{1, \dots, m\}\}$ find quality functions \mathbf{q} , \mathbf{g} of the max-sum classifier such that

$$\mathbf{y}^i = h(\mathbf{x}^i; \mathbf{q}, \mathbf{g}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^{\mathcal{V}}} \left[\sum_{v \in \mathcal{V}} q_v(x_v^i, y_v) + \sum_{\{v, v'\} \in \mathcal{E}} g_{vv'}(y_v, y_{v'}) \right], \quad i \in \mathcal{I}.$$

The max-sum problem $\mathcal{P} = (\mathcal{E}, \mathcal{V}, \mathbf{q}, \mathbf{g}, \mathbf{x})$ associated with the classification $h(\mathbf{x}; \mathbf{q}, \mathbf{g})$ is tractable if:

1. $(\mathcal{V}, \mathcal{E})$ is acyclic graph
2. \mathcal{Y} is fully ordered and $-g_{vv'}$, $\{v, v'\} \in \mathcal{E}$ are submodular w.r.t the ordering: for each $(y_v, y'_v, y_{v'}, y'_{v'}) \in \mathcal{Y}^4$ such that $y_v > y'_v$ and $y_{v'} > y'_{v'}$, it following inequality holds

$$g_{vv'}(y_v, y_{v'}) + g_{vv'}(y'_v, y'_{v'}) \leq g_{vv'}(y_v, y'_{v'}) + g_{vv'}(y'_v, y_{v'})$$

3. $\mathcal{P} = (\mathcal{V}, \mathcal{E}, \mathbf{q}, \mathbf{g}, \mathbf{x})$ have a strictly trivial equivalent, that is, the LP relaxation is tight and the max-sum problem has unique solution