

# Ensemble Methods for Structured Prediction

Vitaly Kuznetsov<sup>1</sup>

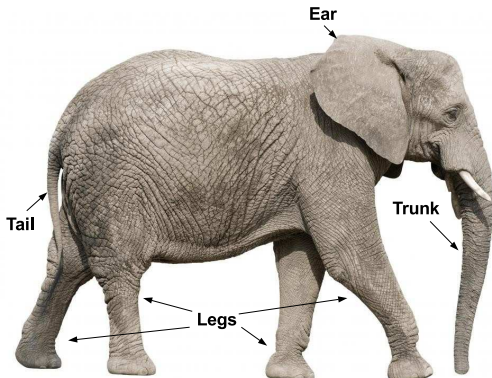
Joint work with

Corinna Cortes<sup>2</sup> and Mehryar Mohri<sup>1,2</sup>

<sup>1</sup>Courant Institute of Mathematical Sciences, New York University

<sup>2</sup>Google Research, New York

# Structured Prediction



# Structured Prediction

Graphemes-to-phonemes task:

- **Input:** sequence of graphemes, e.g.

ensemble algorithm

- **Output:** sequence of phonemes, e.g.

än-'säm-bəl 'al-gə-ri-thəm

# Ensemble Methods

- Often significantly improve performance.
- Benefit from favorable learning guarantees.
- Developed primarily for classification and regression tasks.

# Ensembles & Structured Prediction

Input: ensemble algorithm

Expert 1: än-'səm-bəl 'al-gō-ri-thəm

Expert 2: ən-'säm-bəl əl-gə-ri-thəm

**Goal:** learn to **patch together** predictions of different experts.

# Outline

- Learning scenario.
- Prior work.
- Boosting algorithm.
- On-line solution.
- Experiments.

# Learning Scenario

- Learner receives a sample  $(\mathbf{x}_i, \mathbf{y}_i)_{i=1}^m \in \mathcal{X} \times \mathcal{Y}$ .
- $\mathbf{y} \in \mathcal{Y}$  decomposes into  $\mathbf{y} = (y^1, \dots, y^l)$ .
- Loss is additive

$$L(\mathbf{y}, \tilde{\mathbf{y}}) = \sum_{k=1}^l \ell(y^k, \tilde{y}^k).$$

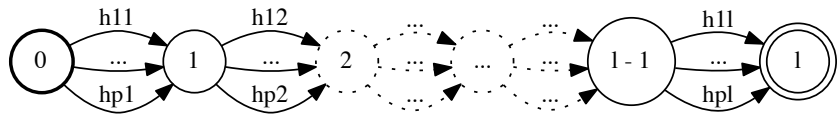
- Learner has access to **black box** predictors  $h_1, \dots, h_p : \mathcal{X} \rightarrow \mathcal{Y}$ .

# Prior Work

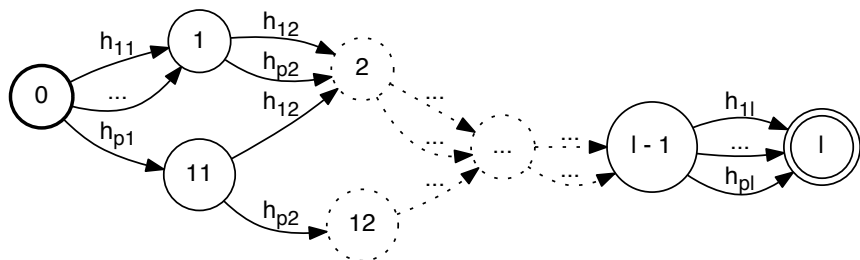
- **Re-ranking techniques:** Collins & Koo, 2005; Huang, 2008.
- **Combinations:** Zeman & Žabokrtský, 2005; Sagae & Lavie, 2006.
- **Scores:** Mohri et. al., 2008; Petrov, 2010; Zhang et. al., 2009.
- **Special experts:** Kocev et. al., 2013; Wang et. al., 2007; Fiscus, 1997.
- **SLE algorithm:** Nguyen & Guo, 2007.



# Path Experts.



# General Graphs



# Boosting Framework

- Learn a scoring function  $\tilde{h} = \sum_{t=1}^T \alpha_t \tilde{h}_t$ ,  $\alpha_t > 0$ .
- Predict

$$\mathcal{H}_{\text{Boost}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \tilde{h}(\mathbf{x}, \mathbf{y}).$$

- No restrictions on base scoring functions  $h_j$ .
- For black box experts

$$\tilde{h}_t(\mathbf{x}, \mathbf{y}) = \sum_{k=1}^l \mathbf{1}_{h_t^k(\mathbf{x})=y^k}.$$

# ESPBoost

**Inputs:** Sample  $S$ ; experts  $\{h_1, \dots, h_p\}$ .

**for**  $i = 1$  **to**  $m$  **and**  $k = 1$  **to**  $l$  **do**

$$\mathcal{D}_1(i, k) \leftarrow \frac{1}{ml}$$

**end for**

**for**  $t = 1$  **to**  $T$  **do**

$$h_t \leftarrow \operatorname{argmin}_{h \in H} \mathbb{E}_{(i,k) \sim \mathcal{D}_t} [\mathbf{1}_{h^k(\mathbf{x}_i) \neq y_i^k}]$$

$$\epsilon_t \leftarrow \mathbb{E}_{(i,k) \sim \mathcal{D}_t} [\mathbf{1}_{h_t^k(\mathbf{x}_i) \neq y_i^k}]$$

$$\alpha_t \leftarrow \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

$$Z_t \leftarrow 2 \sqrt{\epsilon_t (1 - \epsilon_t)}$$

**for**  $i = 1$  **to**  $m$  **and**  $k = 1$  **to**  $l$  **do**

$$\mathcal{D}_{t+1}(i, k) \leftarrow \frac{\exp(-\alpha_t \rho(h_t^k, \mathbf{x}_i, y_i)) \mathcal{D}_t(i, k)}{Z_t}$$

**end for**

**end for**

# ESPBoost algorithm

- Upper bound on empirical loss:

$$\begin{aligned} \frac{1}{ml} \sum_{i=1}^m \sum_{k=1}^l \mathbf{1}_{\mathcal{H}_{\text{Boost}}^k(\mathbf{x}_i) \neq y_i^k} \\ \leq \frac{1}{ml} \sum_{i=1}^m \sum_{k=1}^l \exp \left( - \sum_{t=1}^T \alpha_t \rho(\tilde{\mathbf{h}}_t^k, \mathbf{x}_i, \mathbf{y}_i) \right), \end{aligned}$$

- $\rho(\tilde{\mathbf{h}}_t^k, \mathbf{x}_i, \mathbf{y}_i)$  is margin of the  $\tilde{\mathbf{h}}_t$  at position  $k$  on example  $(\mathbf{x}_i, \mathbf{y}_i)$ .
- ESPBoost algorithm is an application of the coordinate descent to this bound.

# Learning guarantees

## Theorem

Fix  $\rho > 0$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for all  $\mathcal{H}_{\text{Boost}} \in \mathcal{F}$ :

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [L_{\text{Ham}}(\mathcal{H}_{\text{Boost}}(\mathbf{x}), \mathbf{y})] \leq \widehat{R}_\rho \left( \frac{\widetilde{\mathbf{h}}}{\|\boldsymbol{\alpha}\|_1} \right) + \frac{2}{\rho l} \sum_{k=1}^l |\mathcal{Y}_k| \mathfrak{R}_m(H^k) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}},$$

where  $\mathfrak{R}_m(H^k)$  denotes the Rademacher complexity of the class of functions

$$\{\mathbf{x} \mapsto h_j(\mathbf{x}, y) : j \in [1, \rho], y \in \mathcal{Y}_k\}.$$

# Learning guarantees

## Theorem

Let  $\tilde{h}$  denote the scoring function returned by ESPBoost after  $T \geq 1$  rounds. Then, for any  $\rho > 0$ , the following inequality holds

$$\hat{R}_\rho\left(\frac{\tilde{h}}{\|\alpha\|_1}\right) \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\rho} (1 - \epsilon_t)^{1+\rho}}.$$

# On-line Algorithm

Two stage procedure:

1. Run on-line learning algorithm to learn a distribution  $\mathbf{p}$  over path experts.
2. Convert on-line solution  $\mathbf{p}$  to a batch predictor.

Options for on-line algorithms:

Follow-the-Perturbed-Leader (FPL).

Randomized Weighted Majority (RWM).



# RWM Algorithm (Littlestone & Warmuth, 1994)

- $\mathbf{p}_0$  is a uniform distribution over paths experts.
- Receive  $(\mathbf{x}_t, \mathbf{y}_t)$ .
- Update

$$\mathbf{p}_{t+1}(\mathbf{h}) = \frac{\mathbf{p}_t(\mathbf{h})\beta^{L(\mathbf{h}(\mathbf{x}_t), \mathbf{y}_t)}}{\sum_{\mathbf{h}'} \mathbf{p}_t(\mathbf{h}')\beta^{L(\mathbf{h}'(\mathbf{x}_t), \mathbf{y}_t)}}.$$

- Efficient updates using structure of the problem (Takimoto & Warmuth, 2003).

# On-line-to-batch Conversion

- Choose ensemble  $\mathcal{P}$  to minimize:

$$\Gamma(\mathcal{P}) = \frac{1}{|\mathcal{P}|} \sum_{\mathbf{p}_t \in \mathcal{P}} \mathbb{E}_{h \sim p_t} [L(h(x_t), \mathbf{y}_t)] + M \sqrt{\frac{\log \frac{1}{\delta}}{|\mathcal{P}|}}.$$

- Form ensemble distribution  $\mathbf{p} = \frac{1}{|\mathcal{P}|} \sum_{\mathbf{p}_t \in \mathcal{P}} \mathbf{p}_t$ .
- Stochastic or voting predictions.

# Learning Guarantees

## Theorem

For any  $\delta > 0$ , with probability at least  $1 - \delta$  over the choice of the sample  $((\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T))$  drawn i.i.d. according to  $\mathcal{D}$ , the following inequalities hold:

$$\mathbb{E}[L(\mathcal{H}_{\text{Rand}}(\mathbf{x}), \mathbf{y})] \leq \inf_{h \in \mathcal{H}} \mathbb{E}[L(h(\mathbf{x}), \mathbf{y})] + 2M \sqrt{\frac{l \log p}{T}} + 2M \sqrt{\frac{\log \frac{2}{\delta}}{T}}.$$

# Learning Guarantees

## Theorem

*The following inequality relates the generalization error of the majority-vote algorithm to that of the randomized one:*

$$\mathbb{E}[L_{\text{Ham}}(\mathcal{H}_{\text{MVote}}(\mathbf{x}), \mathbf{y})] \leq 2 \mathbb{E}[L_{\text{Ham}}(\mathcal{H}_{\text{Rand}}(\mathbf{x}), \mathbf{y})]$$

*where the expectations are taken over  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$  and  $h \sim p$ .*

# Experiments

Table: Average Normalized Hamming Loss on synthetic data.

	ADS1, $m = 200$	ADS2, $m = 200$
$\mathcal{H}_{\text{MVote}}$	<b>0.0197</b> $\pm$ <b>0.00002</b>	<b>0.2172</b> $\pm$ <b>0.00983</b>
$\mathcal{H}_{\text{Boost}}$	<b>0.0197</b> $\pm$ <b>0.00002</b>	<b>0.2267</b> $\pm$ <b>0.00834</b>
$\mathcal{H}_{\text{SLE}}$	0.5641 $\pm$ 0.00044	0.2500 $\pm$ 0.05003
$\mathcal{H}_{\text{Rand}}$	0.1112 $\pm$ 0.00540	0.4000 $\pm$ 0.00018
BEST $h_j$	0.5635 $\pm$ 0.00004	0.4000

Table: Average Normalized Hamming Loss for ADS3.

$\mathcal{H}_{\text{MVote}}$	<b>0.1788</b> $\pm$ <b>0.00004</b>
$\mathcal{H}_{\text{Boost}}$	0.1831 $\pm$ 0.00240
$\mathcal{H}_{\text{SLE}}$	0.1954 $\pm$ 0.00185
$\mathcal{H}_{\text{Rand}}$	0.3196 $\pm$ 0.00018
Best $h_j$	0.2957 $\pm$ 0.00005

# Experiments

Table: Average Normalized Hamming Loss, Penn Tree Bank.

	TR1, $m = 800$	TR2, $m = 1000$
$\mathcal{H}_{\text{MVote}}$	<b>0.0850</b> $\pm$ <b>0.00096</b>	<b>0.0746</b> $\pm$ <b>0.00014</b>
$\mathcal{H}_{\text{Boost}}$	0.1041 $\pm$ 0.00056	0.1414 $\pm$ 0.00233
$\mathcal{H}_{\text{SLE}}$	<b>0.0778</b> $\pm$ <b>0.00934</b>	<b>0.0814</b> $\pm$ <b>0.02558</b>
$\mathcal{H}_{\text{Rand}}$	0.1128 $\pm$ 0.00048	0.1652 $\pm$ 0.00077
BEST $h_j$	0.1032 $\pm$ 0.00007	0.1415 $\pm$ 0.00005

Table: Average Normalized Hamming Loss for OCR.

$\mathcal{H}_{\text{MVote}}$	0.1992 $\pm$ 0.00274
$\mathcal{H}_{\text{ESPBoost}}$	0.1992 $\pm$ 0.00274
$\mathcal{H}_{\text{SLE}}$	0.1994 $\pm$ 0.00307
$\mathcal{H}_{\text{Rand}}$	0.1994 $\pm$ 0.00276
Best $h_j$	0.1994 $\pm$ 0.00306

# Experiments

Table: Average Normalized Hamming Loss, Pronunciation

	PDS1, $m = 130$	PDS2, $m = 400$
$\mathcal{H}_{\text{MVote}}$	<b>0.2225</b> $\pm$ <b>0.00301</b>	<b>0.2323</b> $\pm$ <b>0.00069</b>
$\mathcal{H}_{\text{Boost}}$	0.3625 $\pm$ 0.01054	0.3499 $\pm$ 0.00509
$\mathcal{H}_{\text{SLE}}$	0.3130 $\pm$ 0.05137	0.3308 $\pm$ 0.03182
$\mathcal{H}_{\text{Rand}}$	0.4713 $\pm$ 0.00360	0.4607 $\pm$ 0.00131
BEST $h_j$	0.3449 $\pm$ 0.00368	0.3413 $\pm$ 0.00067

Table: Average edit-distance, Pronunciation

	PDS1, $m = 130$	PDS2, $m = 400$
$\mathcal{H}_{\text{MVote}}$	<b>0.8395</b> $\pm$ <b>0.01076</b>	<b>0.9626</b> $\pm$ <b>0.00341</b>
$\mathcal{H}_{\text{Boost}}$	1.3977 $\pm$ 0.06017	1.4092 $\pm$ 0.04352
$\mathcal{H}_{\text{SLE}}$	1.1762 $\pm$ 0.12530	1.2477 $\pm$ 0.12267
$\mathcal{H}_{\text{Rand}}$	1.8962 $\pm$ 0.01064	2.0838 $\pm$ 0.00518
BEST $h_j$	1.2163 $\pm$ 0.00619	1.2883 $\pm$ 0.00219

# Experiments

Table: Average Normalized Hamming Loss, Speech

	$p = 5, m = 1500$	$p = 10, m = 1200$
$\mathcal{H}_{\text{MVote}}$	<b>0.2465</b> $\pm$ <b>0.00248</b>	<b>0.2606</b> $\pm$ <b>0.00320</b>
$\mathcal{H}_{\text{Boost}}$	0.2572 $\pm$ 0.00062	0.2864 $\pm$ 0.00103
$\mathcal{H}_{\text{SLE}}$	0.2572 $\pm$ 0.00061	0.2864 $\pm$ 0.00102
$\mathcal{H}_{\text{Rand}}$	0.2877 $\pm$ 0.00480	0.3430 $\pm$ 0.00468
BEST $h_j$	0.2573 $\pm$ 0.00060	0.2865 $\pm$ 0.00101

Table: Average Normalized Hamming Loss, Speech

	$p = 20, m = 900$	$p = 50, m = 700$
$\mathcal{H}_{\text{MVote}}$	<b>0.2773</b> $\pm$ <b>0.00139</b>	<b>0.3217</b> $\pm$ <b>0.00375</b>
$\mathcal{H}_{\text{Boost}}$	0.3115 $\pm$ 0.00089	0.3426 $\pm$ 0.00071
$\mathcal{H}_{\text{SLE}}$	0.3114 $\pm$ 0.00087	0.3425 $\pm$ 0.00076
$\mathcal{H}_{\text{Rand}}$	0.3977 $\pm$ 0.00302	0.4608 $\pm$ 0.00303
BEST $h_j$	0.3116 $\pm$ 0.00087	0.3427 $\pm$ 0.00077



# Non-Additive Losses

- The natural loss functions for most of the key application with structured experts are **non-additive**:
  - machine translation (BLEU score).
  - speech recognition and natural language processing (edit-distance).
  - computational biology ( $n$ -gram similarity measures).
- But existing path experts algorithms cannot be applied with non-additive losses.

# Solution for Non-Additive Losses

- Two new broad families of loss functions: Rational and Tropical Losses.
- Extensions of FPL and RWM to these loss functions based on powerful weighted automata and transducers algorithms.

# Conclusions

- Ensemble methods for structured prediction with learning guarantees.
- On-line and Boosting algorithms.
- Good performance on real and synthetic data.
- Extensions to non-additive losses (e.g. edit-distance).