

Reinforcement Learning for Multi-Step Expert Advice

Patrick Philipp
 Institute AIFB
 Karlsruhe Institute of Technology
 patrick.philipp@kit.edu

Achim Rettinger
 Institute AIFB
 Karlsruhe Institute of Technology
 rettinger@kit.edu

ABSTRACT

Complex tasks for heterogeneous data sources, such as finding and linking named entities in text documents or detecting objects in images, often require multiple steps to be solved in a processing pipeline. In most of the cases, there exist numerous, exchangeable software components for a single step, each an "expert" for data with certain characteristics. Which expert to apply to which observed data instance in which step becomes a challenge that is even hard for humans to decide. In this work, we treat the problem as Single-Agent System (SAS) where a centralized agent learns how to best exploit experts. We therefore define locality-sensitive relational measures for experts and data points, so-called "meta-dependencies", to assess expert performances, and use them for decision-making via Online Model-Free- and Batch Reinforcement Learning (RL) approaches, building on techniques from Contextual Bandits (CBs) and Statistical Relational Learning (SRL). The resulting system automatically learns to pick the best pipeline of experts for a given set of data points. We evaluate our approach for Entity Linking on text corpora with heterogeneous characteristics (such as news articles or tweets). Our empirical results improve the estimation of expert accuracies as well as the out-of-the-box performance of the original experts without manual tuning.

Keywords

Decision-Making with Multi-Step Expert Advice, Expert Processes, Reinforcement Learning, Collective Learning, Entity Linking

1. INTRODUCTION

Complex tasks for heterogeneous data sources often require multiple steps to be solved in a processing pipeline. Entity linking, for example, is a two-step problem where text is tokenized to (i) find named entities and (ii) link them to a knowledge base. In most of the cases, there exist numerous, exchangeable software components for a single step, each an "expert" for data with certain characteristics. We call such a setting *Decision-Making with Multi-Step Expert Advice*.

Combining the outputs of exchangeable experts, e.g. by majority voting, can improve the overall performance, but

Appears in: *Proceedings of the 16th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2017)*, S. Das, E. Durfee, K. Larson, M. Winikoff (eds.), May 8–12, 2017, São Paulo, Brazil.

Copyright © 2017, International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.

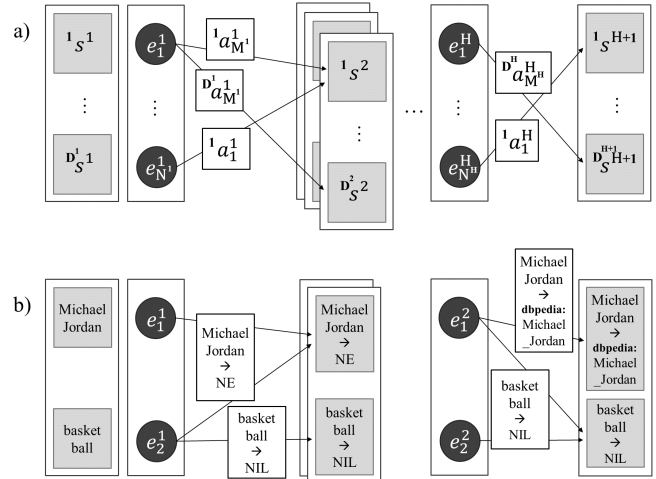


Figure 1: Multi-Step Expert Advice in an Expert Process with a) the formal process and b) an Entity Linking example with two experts.

might fail if experts correlate or display high degrees of variance in their performances. Moreover, two experts solving two subsequent steps might be syntactically compatible while producing bad results when plugged together.

To account for uncertainty, a number of promising candidates might be kept and passed to suitable experts for the subsequent step, instead of using a single output. However, for unstructured data with high dimensionality, such as text or images, the number of potential candidates grows quickly, rendering brute-force approaches impractical for finding optimal configurations. Learning weights for experts is, thus, crucial to query promising experts and keep "good" candidates throughout multiple steps. This also entails learning about impacts of the current decision on future outcomes.

Well-established decision-theoretic frameworks partially deal with these problems, including Decision-Making with Expert Advice (EA) [11] where exchangeable experts are sequentially weighted for single-step problems, Budgeted EA [1] where additional budgets are set on expert queries, Contextual Bandits (CBs) [3] where all experts might be queried but only the subset with the chosen action get feedback and, finally, Contextual Markov Decision Processes (CMDPs) [21] which consider a multi-step CB problem setting, where one outcome is rewarded per step.

Still, none of the frameworks cover *Decision-Making with Multi-Step Expert Advice* at once, leading us to define so-called *Expert Processes* (EPs). EPs enable learning fine-grained expert weights by considering the impact of current decisions on future experts and not being overly restrictive in rewarding outcomes. Based on a budget per step, an agent might query several experts and has to find the optimal joint action, but *observes* all rewards for the queried experts.

To this end, we construct relational measures – so-called “meta-dependencies” – for experts as well as expert pairs, quantifying their predictive performance for specific decision candidates. EPs with meta-dependencies enable the application of techniques from both Reinforcement Learning (RL) as well as Statistical Relational Learning (SRL) to continuously improve decision-making. Relational measures are promising to overcome generalization problems for heterogeneous data distributions, for example available in text or images. These problems occur from assuming data to be independent and identically distributed (i.i.d.), limiting generalization for individual data points. We, therefore, exploit meta-dependencies in an Online- and a Batch RL scenario, where we use the Multiplicative Weights / Hedge (MWH) algorithm [11, 2] for updates and Probabilistic Soft Logic (PSL) [19, 10] for collective inference.

Our main contributions are:

- (i) We provide a formalization for EPs, aligning them with prominent decision-theoretic frameworks.
- (ii) We define and use so-called meta-dependencies as features to assess experts.
- (iii) We present two Model-Free RL approaches – Online RL based on the MWH algorithm and Batch RL using PSL.
- (iv) We instantiate our approaches for Entity Linking and empirically evaluate it against Entity Linking experts as well as other expert combination approaches.

The results clearly show that both our approaches outperform all competitors, but more importantly, are the first to consider meta-dependencies between experts and single data points for learning which expert to pick in a pipeline.

The paper proceeds as follows: We, first, formalize EPs (Sec. 2) and introduce meta-dependencies between experts and data (Sec. 3). We, then, present an Online RL- as well as a Batch RL approach to solve EPs in Sec. 4 and Sec. 5. After describing our experimental scenario and setup, we present and elaborate on our results in Sec. 6. We finally integrate our work into prior research (Sec. 7), summarize it (Sec. 8) and point out possible directions (Sec. 9).

2. EXPERT PROCESSES

An Expert Process EP is a 7-tuple (S, E, A, R, T, H, M) based on a Markov Decision Process (MDP) [31] with S the set of all states, E the set of all experts, A the set of all possible actions, $T : S^h \times A^h \rightarrow S^{h+1}$ the deterministic transition function, $R : S^h \times A^h \rightarrow [0, 1]^{D^h}$ the reward function, H the number of steps and M^h the budget per step. The state set $S = S^1 \cup S^2 \cup \dots \cup S^{H+1}$ consists of $H + 1$ disjunct, heterogeneous state spaces. Hence, a state s^h is uniquely assigned to a step and can be represented in any form, e.g. as a set, sequence, vector, graph or otherwise. Each state s^h

has a set of D^h decision components ${}^1s^h, \dots, {}^{D^h}s^h$, i.e. D^h decisions need to be taken for s^h . As EPs are acyclic and deterministic, choosing $T(s^h, a^h)$ will always result in s^{h+1} . The reward $R(s^h, a^h)$ is a vector of local feedback received after choosing action a^h in state s^h , i.e. the vector contains feedback for each decision component.

We are given a set E of N^h experts, where in step h of the EP a subset of $E = E^1 \cup E^2 \cup \dots \cup E^H$ is available for being queried. An expert e^h is a function $e^h : S^h \rightarrow AR$, where $AR \subseteq S^h \times S^{h+1}$ is an action relation for mapping decision candidates between two state spaces. Possible actions ${}^d a_j^h \in {}^d A^h$ for ${}^d s^h$ are defined as ${}^d A^h := \{{}^d s^{h+1} | \exists e^h \in E : (s^h, s^{h+1}) \in e^h(s^h)\}$. For simplicity, $\Delta_E : S^h \rightarrow S^{h+1}$ returns the suggested action of expert e^h for state s^h , which can be inferred from $e^h(s^h)$. The reward for an expert’s suggested action is, then, $R({}^d s^h, \Delta_{e^h}({}^d s^h))$. The budget M^h confines the number of state-expert pairs (e^h, s^h) one is allowed to query in step h .

The state-action value function (also referred to as Q -function) estimate for taking an action a^h in s^h is defined as:

$$Q(s^h, a^h) = \sum_{{}^d s \in s^h} \sum_{e \in E^h} w_e({}^d s) * \mathbb{1}_{\{\Delta_e({}^d s) = a^h\}} \quad (1)$$

with expert weights:

$$w_{e^h}({}^d s^h) = R({}^d s^h, \Delta_{e^h}({}^d s^h)) + \delta \sum_{a^h \in A_{d_s^h}^h} T(s^{h+1} | s^h, a^h) \max_{a^{h+1} \in A^{h+1}} Q(s^{h+1}, a^{h+1}) \quad (2)$$

where $A_{d_s^h}^h = \{a^h | a^h \in A^h : \Delta_{e^h}({}^d s^h) \in a^h\}$ and δ defines the impact of future rewards. The state-action distribution $P(s^h, a^h)$ is defined as the normalized state-action value function estimate:

$$P(s^h, a^h) = \frac{Q(s^h, a^h)}{\sum_{s' \in \hat{S}^h} \sum_{a' \in A^h} Q(s', a')} \quad (3)$$

where \hat{S}^h is the set of possible states. Similar to MDPs, the Q -function takes into account future actions for the current state-action value. More specifically, learning to balance exploration and exploitation, as prevalent in CBs or MDPs, also appears in EPs in terms of learning which experts to query for which states. For multi-step decision problems, this is referred to as Reinforcement Learning (RL).

Finally, the distribution over (state,expert)-pairs $P(s^h, e^h)$ is constructed based on normalized expert weights:

$$P(s^h, e^h) = \frac{\sum_{{}^d s \in s^h} w_{e^h}({}^d s)}{\sum_{s' \in \hat{S}^h} \sum_{d_{s'} \in s'} \sum_{e' \in E^h} w_{e'}({}^d s')} \quad (4)$$

An EP proceeds in episodes $z \in Z$, where for each episode z the process goes through all steps. In each step h for $h = 1, \dots, H$:

- We are given \hat{S}^h . // dependent on z .
- Build state-expert pair distribution $P(s^h, e^h)$.
- Choose and query M^h state-expert pairs (s^h, e^h) and retrieve a_1^h, \dots, a_M^h .

- Build state-action distribution $P(s^h, a^h)$ based on \hat{S}^h and $\hat{A}^h := \{a^h |^d a^h \in {}^d A^h\}$, containing all possible actions based on a_1^h, \dots, a_M^h .
- Choose (s^h, a^h) and *receive* $R(s^h, a^h)$, but *observe* all residual rewards, i.e. $\forall (\hat{s}^h, \hat{a}^h) : R(\hat{s}^h, \hat{a}^h)$ ¹.
- If $h < H$ then set $\hat{S}^{h+1} = \{s^{h+1} | s^{h+1} \in \hat{A}^h\}$.

The reward of the decision process R^{cum} can be defined as the expected cumulative reward *received* over H , i.e.:

$$R^{cum} = \mathbb{E} \left[\sum_{h=1}^H \sum_{d=1}^{D^h} R^h(d s^h, d a^h) \right] \quad (5)$$

The overall goal of an EP is to learn how to act in the given decision process, i.e. finding the optimal policy π^* with $\pi : S \rightarrow A$ which maximizes R^{cum} .

Our Learning Goal

In this work, we are interested in learning weights w_{e^h} for experts, constrained by a budget M^h . This entails continuously or periodically improving policy π while collecting samples (s^h, e^h) , as the sample collection process directly influences the quality of the learned weights and vice versa. For learning w_{e^h} , we exploit rich state spaces to calculate meta-dependencies between experts and decision candidates.

We will next introduce the concept of meta-dependencies between experts and decision candidates, which is the baseline for learning expert weights for EPs.

3. META-DEPENDENCIES

A meta-dependency expresses a relation between a decision candidate ${}^d s^h$ and either (i) two experts e_i^h, e_j^h of the same step, (ii) two experts e_i^h, e_j^{h+1} of subsequent steps or (iii) a single expert e^h . It approximates how single- or pairs of experts behave for a decision candidate based on behavioural measures for one action $\mu : S^h \times A^h \rightarrow [0, 1]$ or two actions $\mu : S^h \times A^h \times A^h \rightarrow [0, 1]$ to assess ${}^d s_i^h$ on similar decision candidates ${}^c s_j^h$ in its neighbourhood, using a domain dependent kernel $\kappa \in K$ with $\kappa : S^h \times S^h \rightarrow [0, 1]$.

While (iii) is the analogy to function approximation techniques for the meta-analysis of single experts, (i) and (ii) are classes of pairwise measures which enable an agent to decide which experts cooperate well. Given $\mathbf{MD}(e^h, {}^d s^h)$, the set of all calculated meta-dependencies for $(e^h, {}^d s^h)$, we aim to approximate expert weight w_{e^h} assuming different kinds of relationships. We will now present meta-dependencies for each of the three classes.

3.1 Pairwise Intra-Step Experts

Pairwise intra-step meta-dependencies for $(e_i^h, e_j^h, {}^d s^h)$ with $i \neq j$ quantify the relative behaviour between two experts of the same step given a decision candidate, exploiting its neighbourhood. The (ϵ, θ) -dependency of a pair $(e_i^h, e_j^h, {}^d s^h)$ with κ, μ is expressed as:

$$\mathbf{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, {}^d s^h) = \frac{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h) \mu({}^c s^h, \Delta_{e_i^h}({}^c s^h), \Delta_{e_j^h}({}^c s^h))}{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h)} \quad (6)$$

¹Feedback might be available in retrospect at $h = H + 1$.

with $N_{\kappa}({}^d s_i^h) = \{c s_j^h | c s_j^h \in S^h : \kappa({}^d s_i^h, {}^c s_j^h) \geq \gamma \wedge i \neq j\}$ ² the decreasingly sorted neighbourhood of ${}^d s^h$ with minimal similarity γ , $\epsilon = \frac{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h)}{|N_{\kappa}({}^d s^h)|}$ and $|N_{\kappa}({}^d s^h)| \geq \theta$.

To be able to compare two expert actions $\Delta_{e_i^h}({}^d s^h), \Delta_{e_j^h}({}^d s^h)$ we define two behavioural measures for intra-step expert meta-dependencies, joint precision μ_1 and independence μ_2 :

$$\mu_1({}^d s^h, {}^d a_i^h, {}^d a_j^h) = \frac{R({}^d s^h, {}^d a_i^h) + R({}^d s^h, {}^d a_j^h)}{2} \quad (7)$$

$$\mu_2({}^d s^h, {}^d a_i^h, {}^d a_j^h) = |R({}^d s^h, {}^d a_i^h) - R({}^d s^h, {}^d a_j^h)| \quad (8)$$

3.2 Pairwise Inter-Step Experts

Pairwise inter-step meta-dependencies for $(e_i^h, e_j^{h+1}, {}^d s^h)$ with $i \neq j$ quantify the impact of e_i^h on the subsequent step³. The (ϵ, θ) -dependency of a pair $(e_i^h, e_j^{h+1}, {}^d s^h)$ with κ, μ is expressed as:

$$\mathbf{MD}_2^{\mu, \kappa}(e_i^h, e_j^{h+1}, {}^d s^h) = \frac{\sum_{c s^h \in N({}^d s^h)} \kappa({}^d s^h, {}^c s^h) \mu(\Delta_{e_i^h}({}^c s^h), \Delta_{e_j^{h+1}}(\Delta_{e_i^h}({}^c s^h)))}{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h)} \quad (9)$$

with ϵ and θ as defined as before. For inter-step expert meta-dependencies, we define behavioural measure μ_3 as precision of a single expert:

$$\mu_3({}^d s^h, {}^d a^h) = R({}^d s^h, {}^d a^h) \quad (10)$$

3.3 Single Experts

Single expert meta-dependencies for $(e^h, {}^d s^h)$ quantify the behaviour of e^h for ${}^d s^h$ based on similar decision candidates ${}^c s^h$. The (ϵ, θ) -dependency of $(e^h, {}^d s^h)$ with κ, μ is expressed as:

$$\mathbf{MD}_3^{\mu, \kappa}(e^h, {}^d s^h) = \frac{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h) \mu({}^c s^h, \Delta_{e^h}({}^c s^h))}{\sum_{c s^h \in N_{\kappa}({}^d s^h)} \kappa({}^d s^h, {}^c s^h)} \quad (11)$$

with ϵ and θ as defined as before. We use the precision of a single expert μ_3 for single expert meta-dependencies.

We, finally, define a kernel-dependent density measure $\psi_{\kappa}^{\epsilon, \theta}$ to reflect how certain we are in the value of a meta-dependency in the given neighbourhood:

$$\psi_{\kappa}^{\epsilon, \theta}({}^d s^h) = \begin{cases} \frac{N_{\kappa}({}^d s^h)}{\theta} \epsilon, & \text{if } |N_{\kappa}({}^d s^h)| < \theta \\ \epsilon, & \text{otherwise} \end{cases} \quad (12)$$

For expert e^h , the set of meta-dependencies for ${}^d s^h$ is construed as union over all kernels, i.e. $\mathbf{MD}(e^h, {}^d s^h) = \cup_{\kappa \in K} \mathbf{MD}^{\kappa}(e^h, {}^d s^h)$. $\mathbf{MD}^{\kappa}(e^h, {}^d s^h)$ consists of all single-

²For simplicity, we assume that $N_{\kappa}({}^d s_i^h)$ only returns candidates which can be evaluated for a given μ , as numerous state-expert pairs are omitted due to budgets M^h .

³As the current state might contain wrong decisions of prior steps, one could also incorporate the impact of e_j^{h-1} on e_i^h if $\exists c s^{h-1} | e_j^{h-1}({}^c s^{h-1}) = {}^d s^h$.

and pairwise meta-dependencies for κ and the respective behavioural measures μ , instantiated for all residual experts:

$$\begin{aligned} \mathbf{MD}^\kappa(e_i^h, d^h) &= \cup_{\forall e_j^h \in E^h \setminus e_i^h \forall \mu \in \{\mu_1, \mu_2\}} \mathbf{MD}_1^{\mu, \kappa}(e_i^h, e_j^h, d^h) \\ &\cup_{\forall e_j^{h+1} \in E^{h+1}} \mathbf{MD}_2^{\mu_3, \kappa}(e_i^h, e_j^{h+1}, d^h) \\ &\mathbf{MD}_3^{\mu_3, \kappa}(e_i^h, d^h) \end{aligned} \quad (13)$$

To approximate w_{e^h} , single expert meta-dependencies measure the likelihood of an expert's success, while pairwise intra-step expert meta-dependencies use correlations between two experts. As pairwise inter-step expert meta-dependencies take into account future rewards based on experts of the subsequent step, they enable learning Q .

Expert weights can be learned in different RL scenarios. In general, there are Model-based, Value Function-based and Policy Search-based RL approaches, whereas the latter two are Model-Free. Both of our approaches are Value Function-based, as we learn the Q -function based on expert weights. To this end there are two classes of methods, namely Online and Batch RL: Online RL updates Q after each new state, while in Batch RL updates are deferred for a fixed number of episodes to learn supervised models. We next present an Online Model-Free RL approach to learn expert weights.

4. ONLINE RL FOR EXPERT PROCESSES

In the Online RL setting, one aims to continuously improve predictions with every observed reward. We take a Model-Free RL approach, thus not learning R and T (where T is known in our special case) but directly approximating expert weights. Note that most works in RL do not approach generalization across rich world states [21] (although numerous approaches deal with function approximation and limited generalization [35, 25, 40]).

To this end, we use meta-dependencies as direct estimators for expert performances. Over multiple episodes, we keep and update meta-weights for meta-dependencies to gradually get more accurate predictions.

We will first describe the Multiplicative Weights/Hedge (MWH) algorithm and map it to meta-dependencies. Based on the latter, different measures related to Budgeted EA, CBs and MDPs are integrated to account for imperfect information due to the budget and the influence of the current decision on future ones.

4.1 MWH with Meta-Dependencies

The Multiplicative Weights (MW) update method [11, 2] entails drawing experts from a normalized probability distribution and updating their weights according to this probability. Hedging over experts [16] extends MW by exponential weights. For an EP, the essential MWH update rule for $w_{e^h}(z)$ can be defined as:

$$\begin{aligned} w_{e^h}(z+1) &= \\ w_{e^h}(z) \exp\left(\frac{\sum_{s^h \in \hat{S}^h} \sum_{d^h \in S} R^{\text{scaled}}(d^h, \Delta_{e^h}(d^h))}{\sum_{s^h \in \hat{S}^h} \sum_{d^h \in S} \mathbb{1}_{\{q(e^h, s^h)=1\}}} \eta\right) \end{aligned} \quad (14)$$

where $q: E \times S \rightarrow \{0, 1\}$ denotes if expert e was queried for s , $\eta \in [0, 1]$ is the influence parameter of the update and R^{scaled} returns the rewards rescaled to $[-1, 1]$.

The update rule keeps a global weight for each expert as it does not incorporate any available contextual information. To be able to exploit contextual information, we apply MWH directly to meta-dependencies and refer to them as meta-experts $me \in ME$, where $\forall \mathbf{MD}_i(d^h) \in \mathbf{MD}(e^h, d^h) \exists me_i: me_i(d^h) = \mathbf{MD}_i(d^h) \psi_{\kappa}^{\epsilon, \theta}(d^h)$ and $\chi: ME \rightarrow E$ is the assignment function from meta-experts to experts. A meta-expert, thus, predicts the value of meta-dependency $\mathbf{MD}_i(d^h)$ ⁴ for candidate d^h weighted by $\psi_{\kappa}^{\epsilon, \theta}$. Let $R^{\text{meta}}(d^h, me) = me(d^h) R^{\text{scaled}}(d^h, \Delta_{\chi(me)}(d^h))$ be the reward function for meta-experts. The resulting update rule for meta-weights $w_{me}(z+1)$ is:

$$w_{me}(z+1) = w_{me}(z) \exp\left(R^{\text{avg}}(\hat{S}^h, me) \eta\right) \quad (15)$$

where $R^{\text{avg}}(\hat{S}^h, me) = \frac{\sum_{s^h \in \hat{S}^h} \sum_{d^h \in S} R^{\text{meta}}(d^h, me)}{\sum_{s^h \in \hat{S}^h} \sum_{d^h \in S} \mathbb{1}_{\{q(\chi(me), s^h)=1\}}}$. An expert weight, then, is defined by:

$$w_{e^h}(d^h) = \frac{\sum_{me} \mathbb{1}_{\{\chi(me)=e^h\}} me(d^h) w_{me}(z)}{\sum_{me} w_{me}(z)} \quad (16)$$

However, MWH does not deal with incomplete information resulting from budget M^h . We, thus, adapt techniques from Budgeted EA [1] and adversarial CBs to account for state-expert pairs we did not choose to query.

4.2 MWH with Incomplete Information

As an EP allows to query M^h different state-expert pairs, we extend the weight updates with so-called importance weighting, which was proposed for Budgeted EA. This is achieved by normalizing the reward by the probability of all experts queried in z, h with importance weight $iw(z, h) = \sum_{s \in \hat{S}^h} \sum_{e \in E^h} \mathbb{1}_{\{q(e, s)=1\}} P(s, e)$:

$$w_{me}(z+1) = w_{me}(z) \exp\left(\frac{R^{\text{avg}}(\hat{S}^h, me)}{iw(z, h)} \eta\right) \quad (17)$$

Adversarial CBs do not consider data to be i.i.d. but try to gradually adapt their strategies online. The EXP family of CBs (e.g. [3]) are based on MWH, but extend the approach to update expert weights for non-taken actions, i.e. experts in EPs. While in adversarial CBs all experts are queried to suggest an action, all meta-experts are queried in EPs to suggest a weight for their assigned expert.

EXP4.P [7] is an adversarial bandit approach with low regret (i.e. high reward) bound. Here, the probability of choosing an action is equivalent with a normalized expert weight in EPs, which we adapt to:

$$\begin{aligned} w_{e^h}(d^h) &= (1 - N^h p_{\min}^h) \\ &\frac{\sum_{me} \mathbb{1}_{\{\chi(me)=e^h\}} me(d^h) w_{me}(z)}{\sum_{me} w_{me}(z)} + p_{\min}^h \end{aligned} \quad (18)$$

where $p_{\min}^h = [0, \frac{1}{N^h}]$ is the minimum probability for choosing any expert.

The update rule of EXP4.P is a modification of MWH (cmp. eq. 15) and Budgeted EA (cmp. eq. 17) in that

⁴For simplicity $\mathbf{MD}_i(d^h)$ refers to either $\mathbf{MD}_1^{\kappa, \mu}(e_j^h, e_k^h, d^h)$, $\mathbf{MD}_2^{\kappa, \mu}(e_j^h, e_k^{h+1}, d^h)$ or $\mathbf{MD}_3^{\kappa, \mu}(e^h, d^h)$

one ensures exploration for unlikely experts through p_{min}^h and uses confidence bounds to express the variance of the reward. It can be directly used for EPs:

$$w_{me}(z+1) = w_{me}(z) \exp\left(\frac{p_{min}^h}{2} \left(\mathbb{1}_{\{\chi(me)=e\}} \frac{R^{avg}(\hat{S}^h, me)}{iw(z, h)} + \frac{1}{iw(z, h)} \eta \right)\right) \quad (19)$$

Finally, we incorporate outcomes of the chosen joint action by boosting expert weights if an expert's suggested action was wrongfully chosen or wrongfully not chosen, i.e.:

$$\text{wrong}_z^h(me) = \mathbb{1}_{\{q_z^h(\chi(me))=1\}} \left(\mathbb{1}_{\{q_z^h(\Delta_e(s))=0 \wedge (R^{meta}(s, \Delta_e(s)) > 0)\}} \right) + \mathbb{1}_{\{q_z^h(\Delta_e(s))=1 \wedge (R^{meta}(s, \Delta_e(s)) < 0)\}} \quad (20)$$

The resulting boosting factor, then, consists of parameter β which is either active or inactive, i.e. $\text{boost}_z^h(me) = \beta \text{wrong}_z^h(me)$. The final weight update rule is:

$$w_{me}(z+1) = w_{me}(z) \exp\left(\frac{p_{min}^h}{2} \text{boost}_z^h(me) \left(\mathbb{1}_{\{\chi(me)=e\}} \frac{R^{avg}(\hat{S}^h, me)}{iw(z, h)} + \frac{1}{iw(z, h)} \eta \right)\right) \quad (21)$$

Both update rule and expert probability do not specifically incorporate techniques from Model-Free RL approaches, such as the well-known Q -Learning algorithm [41], i.e. future rewards are not explicitly optimized (cmp. eq. 2 for explicit optimization). Here, we rather incorporate future rewards via meta-dependencies, where meta-experts for pairwise inter-step experts take into account how an expert's suggestion influences future results.

Online RL enables reacting to each change in an expert's performance, but often yields suboptimal use of samples. Batch RL [23], on the other hand, enables more efficient use of samples while stabilizing the learning process by updating Q less frequently. Batch RL approaches, therefore, defer updates for a fixed number of episodes.

We, now, take a Batch RL perspective on EPs and use SRL for function approximation.

5. BATCH RL FOR EXPERT PROCESSES

Other than in Online RL, Batch RL [23] approaches conduct updates after a fixed number of episodes Z_{batch} , thereby enabling the use of batch supervised learning. The central idea is to use a scalable SRL approach – i.e. Probabilistic Soft Logic (PSL) – for meta-dependencies to learn a collective model for experts of same or subsequent steps.

We consider an adapted version of an EP where batch updates and -predictions are possible. The new EP is changed to deal with Z_{batch} samples within a single step h before continuing to $h+1$. To use all available information of the EP, we train two models for each step h – an a priori model before querying any expert and an a posteriori model after

having queried all state-expert pairs⁵. To this end, PSL enables directly integrating kernels into the model, thereby transferring inferred weights to other decision candidates by collective inference.

The section proceeds as follows: we first describe PSL, a template language for generating Hinge-Loss Markov Random Fields. We, then, elaborate on how to model weights w_{eh} to approximate Q .

5.1 Probabilistic Soft Logic and Hinge-Loss Markov Random Fields

A Hinge-Loss Markov Random Field (HL-MRF) [5] is a conditional probabilistic model over continuous random variables. HL-MRFs are defined as probability densities:

$$P(Y|X) \propto \frac{1}{Z_{norm}} \exp\left[-\sum_{j=1}^M \lambda_j \phi_j(Y, X)\right]$$

with Z_{norm} the respective normalization constant, weights λ_j , $\phi_j(Y, X) = [\max\{l_j(Y, X), 0\}]^{p_j}$ the hinge-loss potential functions, linear function l_j and $p_j \in \{1, 2\}$.

PSL [19, 10] is a modeling language for HL-MRFs. A model in PSL comprises weighted, first-order logic rules (templates for ϕ_j) with features defining a Markov network. Here, PSL relaxes conjunctions, disjunctions and negations as $A \wedge B = \max\{A+B-1, 0\}$, $A \vee B = \min\{A+B, 1\}$ and $\neg A = 1 - A$ with $A \rightarrow [0, 1]$ and $B \rightarrow [0, 1]$.

The use of hinge-loss feature functions makes inference tractable, as it is a convex optimization problem and thus potentially scales to large datasets. We now present our PSL model based on meta-dependencies for EPs with rich state spaces.

5.2 Meta-Dependencies with PSL

We aim to learn expert weights w_{eh} by using PSL to construct different potential functions $\phi_j(X, Y)$ and learn their weights λ_j . We therefore construct one PSL model per step h and, other than in our online approach, learn one expert weight function for E^h , as our intuition is that meta-dependencies enable to generalize across different experts of the same step. Labels for expert weights are manually constructed based on eq. 2 and normalized *after* all samples are gathered for a batch episode.

All n -ary relations we define map their variables to a real number between zero and one, e.g. Relation : $\text{Var}_i \times \text{Var}_j \rightarrow [0, 1]$ for $n = 2$. An expert weight w_{eh} is represented as relation $W(e^h, d^s)$. Meta-dependency relations for pairwise intra-step experts $\text{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, d^s)$, pairwise inter-step experts $\text{MD}_2^{\mu, \kappa}(e_i^h, e_j^{h+1}, d^s)$ and single experts $\text{MD}_3^{\mu, \kappa}(e^h, d^s)$ are pre-computed meta-dependencies with respective densities $\text{MD}_1 - \text{D}^{\kappa, \mu}(e_i^h, e_j^h, d^s)$, $\text{MD}_2 - \text{D}^{\mu, \kappa}(e_i^h, e_j^{h+1}, d^s)$ and $\text{MD}_3 - \text{D}^{\mu, \kappa}(e^h, d^s)$. Densities are calculated based on $\psi_{\kappa}^{\epsilon, \theta}$ (cmp. eq. 12). $\text{Tr}(e^h, d^s, c^s, c^{s+1})$, the transition relation, denotes that an expert suggests c^s in state d^s .

Pairwise intra-step expert rules express that the weight of an expert increases if (e_i^h, e_j^h, d^s) is (ϵ, θ) -dependent and (e_i^h, e_j^h) agree on d^s with $\mu \in \{\mu_1, \mu_2\}$. We, first, model rules for the a priori setting where experts have not been

⁵Note that we gradually increase our training set with each batch, where densities provide a natural way to passively forget older samples. An alternative would be to actively forget all old batches but reuse old weights as priors.

queried yet. Here, we define relation $\text{Agree}(e_i^h, e_j^h, d^h) = \text{MD}_2^{\kappa, \mu_4}(e_i^h, e_j^h, d^h)$ with $\mu_4(d^h, a_i^h, a_j^h) = \mathbb{1}_{\{d a_i^h = d a_j^h\}}$.

$$\begin{aligned} & \text{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, d^h) \wedge \text{MD}_1 - \text{D}^{\kappa, \mu}(e_i^h, e_j^h, d^h) \\ & \quad \wedge \text{Agree}(e_i^h, e_j^h, d^h) \implies \text{W}(e_i^h, d^h) \\ -\text{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, d^h) \wedge \text{MD}_1 - \text{D}^{\kappa, \mu}(e_i^h, e_j^h, d^h) \\ & \quad \wedge \text{Agree}(e_i^h, e_j^h, d^h) \implies -\text{W}(e_i^h, d^h) \end{aligned}$$

After having queried M^h state-expert pairs, we can leverage new information by only increasing expert weights if there is agreement.

$$\begin{aligned} & \text{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, d^h) \wedge \text{MD}_1 - \text{D}^{\kappa, \mu}(e_i^h, e_j^h, d^h) \\ & \quad \wedge \text{Tr}(e_i^h, d^h, c^h) \wedge \text{Tr}(e_j^h, d^h, c^h) \implies \text{W}(e_i^h, d^h) \\ -\text{MD}_1^{\kappa, \mu}(e_i^h, e_j^h, d^h) \wedge \text{MD}_1 - \text{D}^{\kappa, \mu}(e_i^h, e_j^h, d^h) \\ & \quad \wedge \text{Tr}(e_i^h, d^h, c^h) \wedge \text{Tr}(e_j^h, d^h, c^h) \implies -\text{W}(e_i^h, d^h) \end{aligned}$$

Rules for pairwise inter-step expert meta-dependencies are direct applications of (ϵ, θ) -dependency for (e_i^h, e_j^{h+1}, d^h) and its negative implication with $\mu = \mu_3$. Here, we use the same rules for both the a priori and the a posteriori setting.

$$\begin{aligned} & \text{MD}_2^{\kappa, \mu}(e_i^h, e_j^{h+1}, d^h) \wedge \text{MD}_2 - \text{D}^{\kappa, \mu}(e_i^h, e_j^{h+1}, d^h) \\ & \quad \implies \text{W}(e_i^h, d^h) \\ -\text{MD}_2^{\kappa, \mu}(e_i^h, e_j^{h+1}, d^h) \wedge \text{MD}_2 - \text{D}^{\kappa, \mu}(e_i^h, e_j^{h+1}, d^h) \\ & \quad \implies -\text{W}(e_i^h, d^h) \end{aligned}$$

Similarly, we directly apply (ϵ, θ) -dependency for (e^h, d^h) as positive influence on w_{e^h} with $\mu = \mu_3$, while the negation also holds. A priori and a posteriori models, again, consist of the same rules.

$$\begin{aligned} & \text{MD}_3^{\kappa, \mu}(e^h, d^h) \wedge \text{MD}_3 - \text{D}^{\kappa, \mu}(e^h, d^h) \implies \text{W}(e^h, d^h) \\ -\text{MD}_3^{\kappa, \mu}(e^h, d^h) \wedge \text{MD}_3 - \text{D}^{\kappa, \mu}(e^h, d^h) \implies -\text{W}(e^h, d^h) \end{aligned}$$

We propagate inferred weights by relation $\text{Similar}_\kappa(d^h, c^h)$, thus using kernels directly to exploit similarities among decision candidates. It becomes clear that batch predictions might potentially increase the performance of the approach by clustering similar decision candidates through kernels.

$$\begin{aligned} & \text{Similar}_\kappa(d^h, c^h) \wedge \text{W}(e^h, d^h) \implies \text{W}(e^h, c^h) \\ \text{Similar}_\kappa(d^h, c^h) \wedge -\text{W}(e^h, d^h) \implies -\text{W}(e^h, c^h) \end{aligned}$$

We finally align experts with the same outputs by driving their weights to be the similar. For the a priori setting, we define rules based on the agreement probability:

$$\begin{aligned} & \text{W}(e_i^h, d^h) \wedge \text{Agree}(e_i^h, e_j^h, d^h) \implies \text{W}(e_j^h, d^h) \\ -\text{W}(e_i^h, d^h) \wedge \text{Agree}(e_i^h, e_j^h, d^h) \implies -\text{W}(e_j^h, d^h) \end{aligned}$$

Based on the observed expert action candidates, rules for the a posteriori case are modelled as:

$$\begin{aligned} & \text{W}(e_i^h, d^h) \wedge \text{Tr}(e_i^h, d^h, c^h) \wedge \text{Tr}(e_j^h, d^h, c^h) \\ & \quad \implies \text{W}(e_j^h, d^h) \\ -\text{W}(e_i^h, d^h) \wedge \text{Tr}(e_i^h, d^h, c^h) \wedge \text{Tr}(e_j^h, d^h, c^h) \\ & \quad \implies -\text{W}(e_j^h, d^h) \end{aligned}$$

For learning weights of the PSL model and inferring w_{e^h} , we use an approximate maximum likelihood weight learning algorithm and maximum a posteriori (MAP) inference [4].

6. EMPIRICAL EVALUATION

We evaluate our methods based on an EP for named entity recognition (NER) and disambiguation (NED) (cmp. Fig. 1 b)). The resulting EP has $H = 2$ with S^1 a state space over words, S^2 a state space over named entities and S^3 a state space over disambiguated named entities, i.e. resources available in a (semi-) structured knowledge base (e.g. DBpedia [24], YAGO [36] or Wikipedia⁶). For NER ($h = 1$), decision candidates are n -grams with $n = 1, \dots, |\text{tokens}(s^1)|$. The actual number of decision candidates is determined by E^1 , the set of all NER experts. For NED ($h = 2$), decision candidates are possible named entities, suggested by E^2 , the set of all NED experts. Decision candidates for NED outputs ($h = 3$) are single disambiguated named entities.

6.1 Setup

6.1.1 Data

We cross-validate our approaches on the Microposts 2014 corpus [6] and evaluate on the Spotlight corpus [28] after training on Microposts 2014, where we treat each tweet and each sentence within an article as single state sample. We only use tweets for training phases throughout the evaluation to test the ability of transferring meta-dependencies to different data distributions, as the Spotlight Corpus is comparably small and such problems might often occur in practice.

6.1.2 Meta-Dependencies

Kernels: We summarize the kernel set K and the implementations^{7 8} we used in Tab. 1, where *state* or *candidate* respectively indicates whether the kernel is calculated on a piece of text or a decision candidate.

Kernel	Implementation
Candidate Lingual Type	Part-of-Speech (POS) tag
Candidate Word embeddings	Pretrained on articles ⁷
Candidate MinHash	Jaccard similarity ⁸
State Length	Number of characters
State Extra Characters	Count #, @
State Lingual Type	Average POS tag

Table 1: Used kernels for evaluation.

Densities: We vary θ and γ for the kernels and set:

- $\theta^{\text{length}} = \theta^{\text{extra}} = \theta^{\text{state lingual}} = \frac{Z}{4}$
- $\gamma^{\text{length}} = \gamma^{\text{extra}} = \gamma^{\text{state lingual}} = 0.7$

⁶<http://www.wikipedia.org/>

⁷<https://code.google.com/archive/p/word2vec/>

⁸<https://github.com/ekzhu/datasketch>

- $\theta^{\text{embeddings}} = \theta^{\text{minhash}} = \theta^{\text{candidate lingual}} = \frac{Z}{10}$
- $\gamma^{\text{embeddings}} = \gamma^{\text{minhash}} = \gamma^{\text{candidate lingual}} = 0.5$

6.1.3 Experts

We use expert implementations available as Web services and do not tune parameters.

For NER experts, we use Stanford Tagger [15], FOX [34] and Spotlight Spotter [28]. The NED experts are AGDISTIS [38], AIDA [17] and Spotlight Tagger [28].

6.1.4 Evaluation measures

We evaluate our approach in terms of two target measures: a) accuracy estimation of experts and b) outcome optimization of the EP by performing 10-fold cross-validation and using precision, recall and F1-measure, as proposed by entity linking benchmark GERBIL [39]⁹. For all experiments, we set budgets $M^1 = 2, M^2 = 4$.

a) Accuracy Estimation:

The goal is to correctly estimate the weights w_{e^h} for all experts and, thus, to predict if the suggested action of an expert is correct or not:

- Gl-Avg: Compute global performance weights for experts on left-out training data [33]
- Loc-Avg: Weighted average of single expert meta-dependencies
- On-MWH: Online MWH approach with $p_{min} = 0.05, \eta = 0.8, \beta = 1.5$
- Batch-PSL: Complete PSL model¹⁰ with $Z_{batch} = \frac{Z}{10}$

b) Outcome optimization:

The goal is to optimize the overall outcome of the EP, i.e. R^{cum} :

- AIDA: $e^1 = \text{Stanford Tagger}, e^2 = \text{AIDA}$
- DBS: $e^1 = \text{DBpedia Spotlight Spotter (dictionary lookup)}, e^2 = \text{DBpedia Spotlight Tagger}$
- AGD: $e^1 = \text{FOX}, e^2 = \text{AGDISTIS}$
- Gl-Avg, Loc-Avg, On-MWH, Batch-PSL as defined for a) with majority voting according to $P(s^h, a^h)$.

6.2 Results

We summarize the evaluation results for a) accuracy estimation for $h = 1$ and $h = 2$, and b) outcome optimization in Tab. 2, Tab. 3 and Tab. 4.

The results show that – for a) (accuracy estimation) – our approaches (Batch-PSL and On-MWH) are able to outperform both Gl-Avg without meta-dependencies and Loc-Avg with single expert meta-dependencies. For b) (outcome optimization), Batch-PSL is able to outperform the straightforward use of experts as well as all learning approaches in terms of F1-measure.

⁹For NER a true positive is a correctly found named entity, a true negative is the correctly predicted abstinence of a named entity, a false positive is a mistakenly predicted named entity and a false negative is the mistakenly predicted abstinence of a named entity.

¹⁰<https://github.com/linqs/psl>

Appr.	Microposts '14			Spotlight		
	F1	Prec	Rec	F1	Prec	Rec
Gl-Avg	0.71	0.75	0.67	0.7	0.69	0.72
Loc-Avg	0.76	0.65	0.93	0.78	0.66	0.94
On-MWH	0.85	0.84	0.87	0.82	0.84	0.8
Batch-PSL	0.82	0.8	0.85	0.79	0.78	0.81

Table 2: Evaluation results for a) $h = 1$.

Appr.	Microposts '14			Spotlight		
	F1	Prec	Rec	F1	Prec	Rec
Gl-Avg	0.56	0.63	0.51	0.42	0.39	0.46
Loc-Avg	0.71	0.6	0.87	0.56	0.44	0.78
On-MWH	0.73	0.67	0.81	0.67	0.62	0.72
Batch-PSL	0.8	0.73	0.88	0.74	0.72	0.77

Table 3: Evaluation results for a) $h = 2$.

The results for On-MWH are stable and always close to Batch-PSL, only outperforming the latter for a) (accuracy estimation) for $h = 1$.

In general, our approaches are able to automatically learn fine-grained expert weights and produce good combinations of expert outputs.

6.2.1 Qualitative Results of Meta-Dependencies

The learned meta-weights show that single- as well as pairwise intra-step expert meta-dependencies both drive weights towards one, while pairwise inter-step expert meta-dependencies have the opposite effect. The meta-dependencies, thus, complement each other to provide well-balanced estimates.

Although kernels for state text length and -extra characters are heuristics, they improve the predictive power. As both the PSL model as well as the adapted version of MWH are easily extensible and scale well, adaptation to new tasks or domains is straightforward.

The textual kernels are most influential, suggesting that additional textual kernels, covering more decision candidates, might substantially improve the predictions. Our approaches, thus, enable exploiting numerous embedded text representations to be combined for specific complex tasks.

Finally, the results suggest that meta-dependencies learned on tweets also improve learning expert weights for news articles, making knowledge transfer possible.

7. RELATED WORK

We, first, situate EPs into well-established decision-theoretic frameworks of SAS and compare our approaches to related SAS techniques. We, then, integrate our work into related MAS branches, as we argue that our approaches might lead towards better coordination and communication for solving EPs in decentralized MAS.

7.1 Single-Agent Systems

Decision-Making with Multi-Step Expert Advice directly maps to SAS where one agent controls all actions (i.e. experts) and is related to (i) Decision-Making with Expert Advice (EA) [11], as we are able to observe the rewards of all queried experts. However, since there is a budget for querying experts we need to take decisions inbetween steps. There are, thus, overlaps with (ii) Budgeted EA [1] where

Appr.	Microposts '14			Spotlight		
	F1	Prec	Rec	F1	Prec	Rec
DBS	0.33	0.24	0.51	0.45	0.38	0.57
AGD	0.28	0.59	0.19	0.13	0.56	0.07
AIDA	0.33	0.69	0.21	0.24	0.69	0.15
GI-Avg	0.41	0.51	0.34	0.36	0.48	0.29
Loc-Avg	0.47	0.59	0.38	0.42	0.49	0.37
On-MWH	0.54	0.6	0.5	0.43	0.44	0.42
Batch-PSL	0.56	0.72	0.46	0.48	0.64	0.39

Table 4: Evaluation results for b)

not all experts can be queried and (iii) Contextual Bandits (CBs) [3] which reward only one action. In addition, as we deal with decision processes over multiple steps, (iv) Contextual Markov Decision Processes (CMDPs) [21] are strongly related to EPs, as they consider multi-step CB problems. In comparison, in EPs one is allowed to choose multiple decision candidates to query a subset of experts and needs to keep respective probability distributions. While a single action needs to be chosen in both EPs and CMDPs, one receives rewards for all available actions in EPs.

Knowledge-based trust [13] extends knowledge fusion approaches [12], which combine the outputs of multiple triple extractors from multiple websites and is thus an instance of (i), where a combination function over experts is learned based on heterogeneous data. While we do not deal with problem settings where input data might be flawed, our approach could be integrated into knowledge-based trust by providing accurate confidence measures for available information extractors.

To determine the accuracy of experts in an independent and identically distributed (i.i.d.) data setting, [29, 30] use pairwise error independence measures between binary classifiers without requiring labelled data and also fall under (i). We do not solely rely on unlabeled data, but leverage Collective Learning to integrate labeled and unlabeled data. In addition, experts are allowed to perform worse than average.

The algorithm selection problem has widely been studied for combinatorial search problems [20] or automatic machine learning (AutoML) [37], and was originally stated by [32]. Algorithm selection problem is a special case of (ii) and (iii), and AutoML falls under (iv) as very rich state spaces have to be dealt with, multiple steps partially occur and a budget is imposed during the learning phase. Automating machine learning pipelines, as approach by auto-sklearn [14], entails automating feature selection, algorithm selection and hyperparameter optimization. While prominent Bayesian optimization techniques share numerous overlaps with strategies to balance exploitation and exploration in CBs or MDPs, our focus is guiding the decision process for single states and their decision candidates.

While ensemble learning deals with learning functions over individual models [43], meta-learning aims to optimize the application of learned models to new data sets based on performance histories and so-called meta-features describing the datasets [9]. Both types of approaches can be seen as batch versions of (i). Our work, however, focuses on selecting experts (i.e. algorithms) for specific data points by exploiting meta-dependencies, where multiple instances of the same expert with different parameterizations could be used.

7.2 Multi-Agent Systems

While we defined EPs for centralized decision-making in SAS, a MAS view would distribute control of one or more experts to decentralized agents, rendering agent communication and coordination important. Here, one would deal with a multi-stage MAS where agents are fully cooperative. Works for single-stage coordination [18, 26] in MAS deal with agents for one-step problems, but extensions to the multi-stage case exist [8]. The latter are MDP extensions and partition the action space according to available agents.

Communication might be limited if agents cannot query their experts in each episode due to a budget or if agents have to choose a subset of agents to communicate with (eg. [42]). To this end, central differences occur for payoffs, where in cooperative scenarios either the joint action is rewarded for all agents [18] (i.e. they receive the same payoffs) or, such as in bandit settings, only those agents receive a reward who voted for the chosen action [26].

8. CONCLUSION

We formalized Expert Processes for Multi-Step Expert Advice by reusing central concepts of Decision-Making with Expert Advice, -Budgeted Expert Advice, Markov Decision Processes and Contextual Bandits. For learning distribution-independent expert weights, we introduced meta-dependencies as features, which constitute diverse relations among single- or pairs of experts and data points. We argued that meta-dependencies can be best exploited by Online Model-Free Reinforcement Learning based on the Multiplicative Weights/Hedge algorithm and Batch Model-Free Reinforcement Learning with a Probabilistic Graphical Model as function approximator, both enabling to learn expert weights suitable for Expert Processes. We evaluated our approaches for an Entity Linking Expert Process, where they were able to provide superior results for estimating accuracies of experts and optimizing the overall outcome for tweets and news articles.

9. FUTURE WORK

EPs enable applying RL approaches to Multi-Step Expert Advice. As our Online Model-Free RL approach is based on EA, we do not directly deal with function approximation. There are extensions to well-known RL approaches, such as Q-learning with linear function approximation [27] or Least Squares Policy Iteration (LSPI) [22, 25], which might be helpful to learn functions over meta-experts.

We solve EPs in a SAS, assuming a central agent to decide which experts to query for which states. However, a decentralized MAS might enable higher degrees of scalability due to parallelization and higher degrees of generalizability of learned weights due to additional self-interest of each agent. Approaches for learning decentralized coordination [42] or learning to communicate in MAS [8, 18] might, thus, be beneficial for EPs. In addition, pairwise meta-dependencies, e.g. for intra- or inter-step experts, seem to be readily applicable for optimizing coordination and communication in such a setting.

REFERENCES

- [1] K. Amin, S. Kale, G. Tesauro, and D. S. Turaga. Budgeted prediction with expert advice. In *AAAI*, pages 2490–2496, 2015.

- [2] S. Arora, E. Hazan, and S. Kale. The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(1):121–164, 2012.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire. The non-stochastic multiarmed bandit problem. *SIAM J. Comput.*, 32(1):48–77, 2002.
- [4] S. H. Bach, M. Broecheler, L. Getoor, and D. P. O’Leary. Scaling MPE inference for constrained continuous markov random fields with consensus optimization. In *NIPS*, pages 2663–2671, 2012.
- [5] S. H. Bach, B. Huang, B. London, and L. Getoor. Hinge-loss markov random fields: Convex inference for structured prediction. In *UAI*, 2013.
- [6] A. E. C. Basave, G. Rizzo, A. Varga, M. Rowe, M. Stankovic, and A. Dadzie. Making sense of microposts (#microposts2014) named entity extraction & linking challenge. In *Making Sense of Microposts (WWW)*, pages 54–60, 2014.
- [7] A. Beygelzimer, J. Langford, L. Li, L. Reyzin, and R. E. Schapire. Contextual bandit algorithms with supervised learning guarantees. In *AISTATS*, pages 19–26, 2011.
- [8] C. Boutilier. Sequential optimality and coordination in multiagent systems. In *IJCAI*, pages 478–485, 1999.
- [9] P. Brazdil, C. Giraud-Carrier, C. Soares, and R. Vilalta. *Metalearning: Applications to Data Mining*. Springer, 1 edition, 2008.
- [10] M. Bröcheler, L. Mihalkova, and L. Getoor. Probabilistic similarity logic. In *UAI*, pages 73–82, 2010.
- [11] N. Cesa-Bianchi, Y. Freund, D. Haussler, D. P. Helmbold, R. E. Schapire, and M. K. Warmuth. How to use expert advice. *J. ACM*, 44(3):427–485, 1997.
- [12] X. L. Dong, E. Gabrilovich, G. Heitz, W. Horn, K. Murphy, S. Sun, and W. Zhang. From data fusion to knowledge fusion. *PVLDB*, 7(10):881–892, 2014.
- [13] X. L. Dong, E. Gabrilovich, K. Murphy, V. Dang, W. Horn, C. Lugaresi, S. Sun, and W. Zhang. Knowledge-based trust: Estimating the trustworthiness of web sources. *PVLDB*, 8(9):938–949, 2015.
- [14] M. Feurer, A. Klein, K. Eggenberger, J. T. Springenberg, M. Blum, and F. Hutter. Efficient and robust automated machine learning. In *NIPS*, pages 2962–2970, 2015.
- [15] J. R. Finkel, T. Grenager, and C. D. Manning. Incorporating non-local information into information extraction systems by gibbs sampling. In *ACL*, 2005.
- [16] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119 – 139, 1997.
- [17] J. Hoffart, M. A. Yosef, I. Bordino, H. Fürstenau, M. Pinkal, M. Spaniol, B. Taneva, S. Thater, and G. Weikum. Robust disambiguation of named entities in text. In *EMNLP*, pages 782–792, 2011.
- [18] S. Kapetanakis and D. Kudenko. Reinforcement learning of coordination in cooperative multi-agent systems. In *AAAI*, pages 326–331, 2002.
- [19] A. Kimmig, S. Bach, M. Broecheler, B. Huang, and L. Getoor. A short introduction to probabilistic soft logic. In *NIPS Workshop on Probabilistic Programming: Foundations and Applications*, pages 1–4, 2012.
- [20] L. Kotthoff. Algorithm selection for combinatorial search problems: A survey. *AI Magazine*, 35(3):48–60, 2014.
- [21] A. Krishnamurthy, A. Agarwal, and J. Langford. PAC reinforcement learning with rich observations. In *NIPS*, pages 1840–1848, 2016.
- [22] M. G. Lagoudakis and R. Parr. Least-squares policy iteration. *Journal of Machine Learning Research*, 4:1107–1149, 2003.
- [23] S. Lange, T. Gabel, and M. Riedmiller. *Batch Reinforcement Learning*, pages 45–73. Springer Berlin Heidelberg, 2012.
- [24] J. Lehmann, R. Isele, M. Jakob, A. Jentzsch, D. Kontokostas, P. N. Mendes, S. Hellmann, M. Morsey, P. van Kleef, S. Auer, and C. Bizer. Dbpedia - A large-scale, multilingual knowledge base extracted from wikipedia. *Semantic Web*, 6(2):167–195, 2015.
- [25] L. Li, M. L. Littman, and C. R. Mansley. Online exploration in least-squares policy iteration. In *AAMAS*, pages 733–739, 2009.
- [26] K. Liu and Q. Zhao. Distributed learning in multi-armed bandit with multiple players. *Trans. Sig. Proc.*, 58(11):5667–5681, 2010.
- [27] F. S. Melo and M. I. Ribeiro. Q-learning with linear function approximation. In *COLT*, pages 308–322, 2007.
- [28] P. N. Mendes, M. Jakob, A. García-Silva, and C. Bizer. Dbpedia spotlight: shedding light on the web of documents. In *I-SEMANTICS*, pages 1–8, 2011.
- [29] A. Platanios, A. Blum, and T. M. Mitchell. Estimating accuracy from unlabeled data. In *In Proceedings of UAI*, 2014.
- [30] E. A. Platanios, A. Dubey, and T. M. Mitchell. Estimating accuracy from unlabeled data: A bayesian approach. In *ICML*, pages 1416–1425, 2016.
- [31] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994.
- [32] J. R. Rice. The algorithm selection problem. *Advances in Computers*, 15:65–118, 1976.
- [33] P. Ruiz and T. Poibeau. Combining open source annotators for entity linking through weighted voting. In *SEM*, 2015.
- [34] R. Speck and A. N. Ngomo. Ensemble learning for named entity recognition. In *In ISWC*, pages 519–534, 2014.
- [35] A. L. Strehl, L. Li, E. Wiewiora, J. Langford, and M. L. Littman. PAC model-free reinforcement learning. In *ICML*, pages 881–888, 2006.
- [36] F. M. Suchanek, G. Kasneci, and G. Weikum. Yago: a core of semantic knowledge. In *WWW*, pages 697–706, 2007.
- [37] C. Thornton, F. Hutter, H. H. Hoos, and K. Leyton-Brown. Auto-weka: Combined selection and hyperparameter optimization of classification algorithms. In *SIGKDD*, pages 847–855. ACM, 2013.
- [38] R. Usbeck, A. N. Ngomo, M. Röder, D. Gerber, S. A.

- Coelho, S. Auer, and A. Both. AGDISTIS - graph-based disambiguation of named entities using linked data. In *ISWC*, pages 457–471, 2014.
- [39] R. Usbeck, M. Röder, A. N. Ngomo, C. Baron, A. Both, M. Brümmer, D. Ceccarelli, M. Cornolti, D. Cherix, B. Eickmann, P. Ferragina, C. Lemke, A. Moro, R. Navigli, F. Piccinno, G. Rizzo, H. Sack, R. Speck, R. Troncy, J. Waitelonis, and L. Wesemann. GERBIL: general entity annotator benchmarking framework. In *WWW*, pages 1133–1143, 2015.
- [40] T. J. Walsh, I. Szita, C. Diuk, and M. L. Littman. Exploring compact reinforcement-learning representations with linear regression. In *UAI*, pages 591–598, 2009.
- [41] C. J. C. H. Watkins and P. Dayan. Q-learning. *Machine Learning*, 8(3):279–292, 1992.
- [42] P. Xuan, V. R. Lesser, and S. Zilberstein. Communication decisions in multi-agent cooperation: model and experiments. In *Agents*, pages 616–623, 2001.
- [43] Z.-H. Zhou. *Ensemble methods: foundations and algorithms*. CRC press, 2012.