

Factored Latent-Dynamic Conditional Random Fields for Single and Multi-label Sequence Modeling

Satyajit Neogi

*School of Electrical and Electronic Engineering
Nanyang Technological University
50 Nanyang Avenue, Singapore - 639798*

SATYAJIT001@E.NTU.EDU.SG

Justin Dauwels

*School of Electrical and Electronic Engineering
Nanyang Technological University
50 Nanyang Avenue, Singapore - 639798*

JDAUWELS@NTU.EDU.SG

Editor:

Abstract

Conditional Random Fields (CRF) are frequently applied for labeling and segmenting sequence data. Morency et al. (2007) introduced hidden state variables in a labeled CRF structure in order to model the latent dynamics within class labels, thus improving the labeling performance. Such a model is known as Latent-Dynamic CRF (LDCRF). We present Factored LDCRF (FLDCRF), a structure that allows multiple latent dynamics of the class labels to interact with each other. Including such latent-dynamic interactions leads to improved labeling performance on single-label and multi-label sequence modeling tasks. We apply our FLDCRF models on two single-label (one nested cross-validation) and one multi-label sequence tagging (nested cross-validation) experiments across two different datasets - UCI gesture phase data and UCI opportunity data. FLDCRF outperforms all state-of-the-art sequence models, i.e., CRF, LDCRF, LSTM, LSTM-CRF, Factorial CRF, Coupled CRF and a multi-label LSTM model in all our experiments. In addition, LSTM based models display inconsistent performance across validation and test data, and pose difficulty to select models on validation data during our experiments. FLDCRF offers easier model selection, consistency across validation and test performance and lucid model intuition. FLDCRF is also much faster to train compared to LSTM, even without a GPU. FLDCRF outshines the best LSTM model by $\sim 4\%$ on a single-label task on UCI gesture phase data and outperforms LSTM performance by $\sim 2\%$ on average across nested cross-validation test sets on the multi-label sequence tagging experiment on UCI opportunity data. The idea of FLDCRF can be extended to joint (multi-agent interactions) and heterogeneous (discrete and continuous) state space models.

Keywords: Conditional Random Fields, Sequence Labeling, Multi-task learning, Latent-dynamic models.

1. Introduction

Labeling and segmenting sequence data is a very common problem of machine learning, which has various applications in the fields of Natural Language Processing (e.g., noun phrase chunking (Sutton et al., 2007; Collobert et al., 2011), Part of Speech tagging (Sutton et al., 2007; Collobert et al., 2011), named entity recognition (Lample et al., 2016;

Collobert et al., 2011) etc.), computer vision (gesture recognition (Morency et al., 2007), activity recognition (Ordez and Roggen, 2016) etc.), and information extraction (McCallum et al., 2000). Given an observed input sequence $\mathbf{x} = \{x_t\}_{t=1:T}$, the labeling task aims at automatically assigning labels y_t for each input value x_t .

Conditional Random Fields (CRF, Lafferty et al. (2001)) are frequently applied to the sequence labeling task. The simplest of its kind is the Linear Chain Conditional Random Field (LCCRF, also called CRF), which is defined over $\{x_t\}_{t=1:T}$ and $\{y_t\}_{t=1:T}$ by the graphical constraints depicted in Fig. 1a. The output sequence values (y_t) usually belong to a predefined set \mathcal{Y} of class labels. For example, in a continuous human activity recognition problem, $\mathcal{Y} = \{\text{‘sitting’}, \text{‘standing’}, \text{‘walking’}, \text{‘lying’}\}$. Due to the Markov dependency assumption among the label sequence $\{y_t\}_{t=1:T}$, a LCCRF models the extrinsic dynamics (temporal dependency) within the class labels.

In addition to the extrinsic dynamics among the class labels, each class in \mathcal{Y} has underlying intrinsic dynamics. For instance, each class ‘sitting’, ‘standing’, ‘walking’ etc. contains temporal transitions between certain hidden posture states. Morency et al. (2007) introduced a layer of hidden variables $\{h_t\}_{t=1:T}$ (see Fig. 1c) in a LCCRF structure, in order to model a complete latent (both extrinsic and intrinsic) dynamics of the class labels. Each class in \mathcal{Y} is associated with distinct sets of hidden states. Transitions among the states (to capture the intrinsic and extrinsic dynamics) are made possible through the layer of hidden variables $\{h_t\}_{t=1:T}$. LDCRF has been shown to outperform popular sequence models viz., Hidden Markov Models, LCCRF and Hidden-state CRFs (see Section 2) on several sequence labeling tasks (Morency et al., 2007; Sun et al., 2008). We describe the LDCRF model in Appendix A.

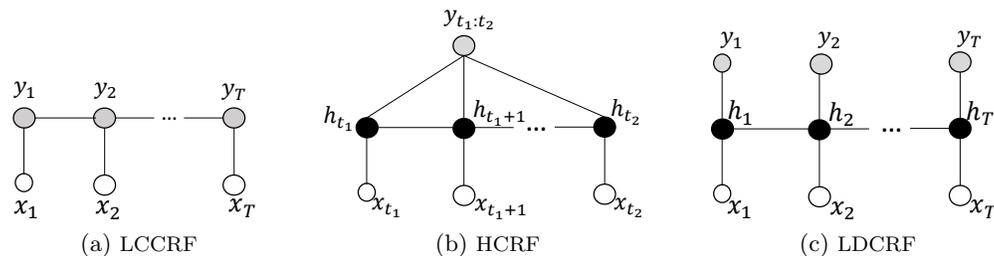


Figure 1: Single-label sequential CRF variants. White nodes are (training + testing) observed and black nodes are hidden. Grey nodes are observed only during training.

Learning performance in a LDCRF can only be improved by varying the number N_s of states associated to each class label. This only mode of variation restricts the model capabilities, and results in a rapidly growing state transition matrix with size $(N_l \cdot N_s) \times (N_l \cdot N_s)$, where $N_l = |\mathcal{Y}|$. Such increment results in greater model complexity and requires more training data. A solution for a Hidden Markov Model with a similar structure of the state space was proposed by Ghahramani and Jordan (1997), where multiple cotemporal state variables ($h_{1,t}$, $h_{2,t}$ etc.) replaced a single hidden state variable (see Fig. 2). The hidden states are distributed among the state variables and all state variables assumed independent temporal dynamics over the input observations $\{x_t\}_{t=1:T}$, resulting in reduction in number of state transition parameters.

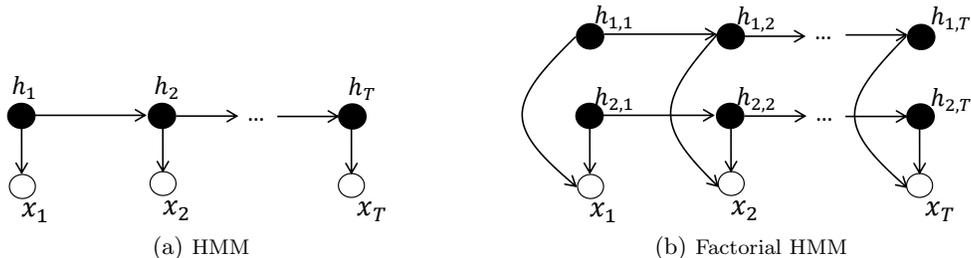


Figure 2: (a) A Hidden Markov Model. (b) A Factorial HMM with multiple cotemporal state variables.

In addition, it is possible that there exists multiple interacting latent dynamics within the class labels. For example, in a continuous human activity recognition problem (e.g., with four classes, ‘lying’, ‘sitting’, ‘standing’ and ‘walking’), a human being has two interacting dynamics - one along the plane and another in vertical direction. Such interaction among different latent dynamics can be captured by connecting state variables ($h_{1,t}$, $h_{2,t}$ etc.) across layers (see Fig. 3b). LDCRF (see Fig. 3a) ignores any possible interaction among the associated hidden states at slice t , and thus is unable to capture such interactions. Considering all these limitations of LDCRF for a single label sequence modeling and inspired by Factorial HMM, we present a generalization of LDCRF, called Factored Latent-Dynamic Conditional Random Fields (FLDCRF, Neogi et al. (2017)), in order to:

- generate new models by varying the number of hidden layers,
- generate factorized models (Ghahramani and Jordan, 1997) with fewer parameters, which need less data for training, and
- model multiple interacting latent dynamics within class labels.

We denote the FLDCRF single-label sequence model (see Fig. 3b) as FLDCRF-s. FLDCRF-s improves LDCRF performance across four nested cross-validation experiments on UCI gesture phase data (Madeo et al., 2013) and UCI opportunity data (Chavarriaga et al., 2013) (see Section 6).

LSTMs (Hochreiter and Schmidhuber, 1997) are the most popular kind of Recurrent Neural Networks (RNN) for sequence modeling, largely due to their ability to capture long-range dependencies among sequence values. They are frequently applied to sequence labeling tasks, viz., named entity recognition (Lample et al., 2016; Huang et al., 2015), noun phrase chunking (Huang et al., 2015), action recognition (Ullah et al., 2017) etc. A variant of LSTM, combined with a CRF layer (Lample et al., 2016) for classification, has recently been very successful in sequence labeling tasks (Lample et al., 2016; Huang et al., 2015; Zhu et al., 2019). We compare FLDCRF-s with both LSTM and LSTM-CRF on the test data in all our experiments. In addition, we analyze several important modeling aspects of FLDCRF-s and LSTM, viz., ease of model selection, consistency across validation and test performance, computation times for training and inference etc. FLDCRF-s not only outperforms LSTM and LSTM-CRF on test data across all experiments, but also offers easier model selection and more consistent performance across validation and test data (see

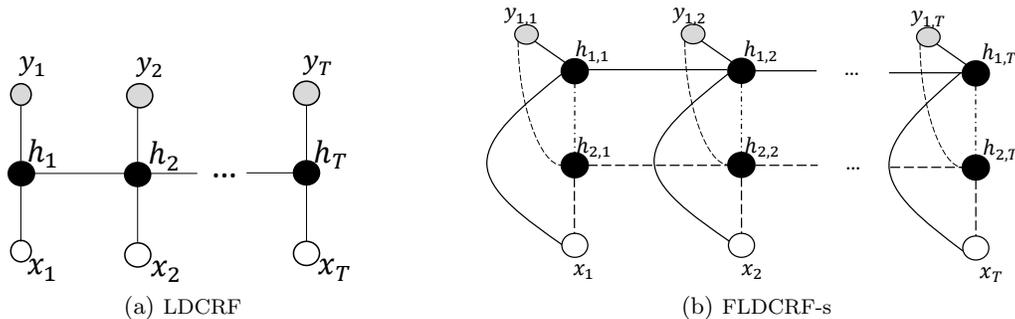


Figure 3: (a) A Latent-Dynamic Conditional Random Field. (b) FLDCRF graphical model for single-label sequence prediction. Solid and dashed connections depict two different latent dynamics within class labels $\{y_t\}_{t=1:T}$. Dashed-dotted inter-layer connections capture interactions among different latent dynamics. The graph shows only two hidden layers and Markov connections for transitions between state variables $\{h_{i,t}\}_{t=1:T}$ in layer i , $i = 1 : L$. More hidden layers and long-range dependencies (semi-Markov for transitions and Markov/semi-Markov for dashed-dotted inter-layer influences) are also possible but omitted for simplicity.

Section 6). In addition, FLDCRF-s requires notably less training and inference times than LSTM models, even without GPU implementation.

Multi-task sequence learning (Sutton et al., 2007; Collobert et al., 2011) is the task to jointly tag sequence values with multiple label categories. Such learning is often helpful where different label categories are related to each other. For example, let us consider the two tasks of continuously recognizing a high-level action (e.g., with two classes, ‘relaxing’ and ‘exercising’) and a low-level action (e.g., with four classes, ‘lying’, ‘sitting’, ‘standing’ and ‘walking’) over some given input (sensor) sequence data. Very clearly the two action types are related and a joint modeling of both is likely to improve the individual recognition performances with the help of additional contextual information (see Section 6).

Dynamic CRFs (see Fig. 4b-c, Sutton et al. (2007)) are very popular (Sutton et al., 2007; Lu et al., 2010) for joint modeling of multiple label sequences. Different label categories $y_{1,t}$, $y_{2,t}$ etc. interact with each other through inter-layer links. The structure in Fig. 4b is popularly known as a Factorial CRF.

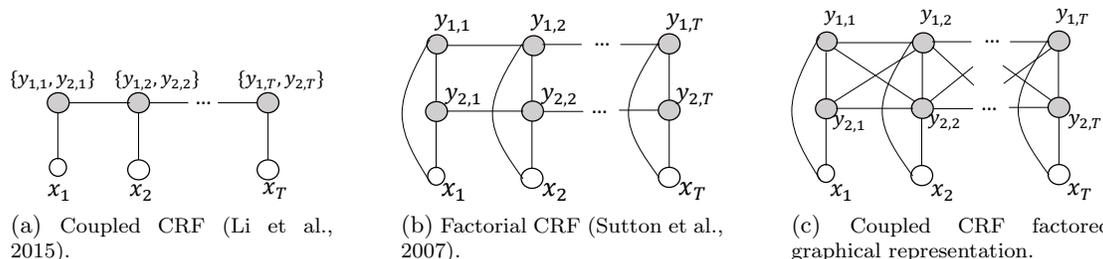


Figure 4: Existing multi-label sequence learning models. (b) and (c) are examples of Dynamic Conditional Random Fields (Sutton et al., 2007).

Similar to a Linear chain CRF, a DCRF does not contain any hidden state variables and therefore is unable to capture the latent (intrinsic and extrinsic) dynamics within the classes

of different label categories (e.g., high and low-level actions). Moreover, connections between the latent dynamics of different label categories can possibly capture deeper understanding of their interactions. Considering such limitation of DCRF, we propose a multi-label variant of FLDCRF (FLDCRF-m, see Fig. 5a), in order to:

- introduce latent variables in a general DCRF structure,
- model interactions among the latent dynamics of two or more inter-related label categories, and
- provide a general representation of DCRFs with embedded hidden state variables for each output variable. The original DCRF structure can be recovered by associating one hidden state per class of the output variables (see Section 3.1 for mathematical details).

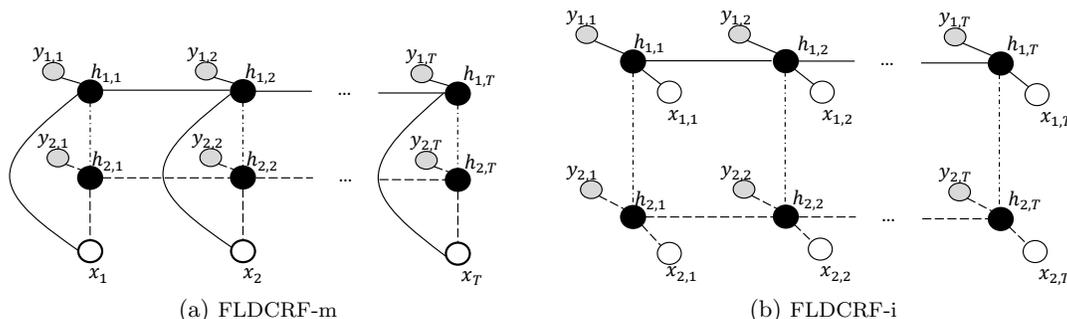


Figure 5: Variants of FLDCRF. (a) FLDCRF graphical model for multi-label sequence prediction. Different label categories, $y_{1,t}$ and $y_{2,t}$, over input x_t are connected through their respective hidden layers, $h_{1,t}$ and $h_{2,t}$, influencing each others’ latent dynamics. (b) FLDCRF graphical model for social interaction. Two different objects, $(\{x_{1,t}\}, \{y_{1,t}\})$ and $(\{x_{2,t}\}, \{y_{2,t}\})$ are shown to be interacting in the latent space through their hidden layers to effect each others’ dynamics. Longer range dependencies are possible in both models but are omitted for simplicity. In fact, we add markov connections for influence links in our multi-label experiment (see Section 5).

A less popular CRF variant for joint sequence tagging is coupled CRF (CCRF, see Fig. 4, Li et al. (2015)), that utilizes a combined state space for different label categories $y_{1,t}$, $y_{2,t}$ etc. Recurrent Neural Networks (Pahuja et al., 2017) have also been applied to joint tagging of sequences. We present a structured and detailed review of literature on joint sequence labeling in Section 2.

FLDCRF-m outperforms Factorial CRF and coupled CRFs on our multi-task sequence labeling experiment on the UCI opportunity dataset (Chavarriaga et al., 2013) (see Section 6). FLDCRF-m also outshines all tested single-label sequence models (trained and inferred separately on different label categories), viz., CRF, LDCRF, FLDCRF-s, LSTM and LSTM-CRF on both tasks combined, while improving the individual labeling performances on most nested test sets. We also compare FLDCRF-m with a multi-task LSTM model, which we refer to as LSTM-m in the paper, similar to the BiRNN model by Pahuja et al. (2017).

Another variant of FLDCRF, called the interaction variant (FLDCRF-i, see Fig. 5b), captures relationship among the dynamics of multiple agents in a social environment.

FLDCRF-i can be applied to multi-agent sequence modeling and prediction tasks in a social environment, e.g., joint intention prediction of multiple pedestrians). Social LSTM (Alahi et al., 2016) provides a similar unsupervised framework for joint path prediction. We present a general mathematical description of all the 3 FLDCRF model variants (FLDCRF-s, FLDCRF-m and FLDCRF-i) in Section 3.

We summarize the main contributions of this paper below:

- We propose a generalization of LDCRF, called Factored LDCRF, to capture sequential interactions in latent space.
- Single label variant of FLDCRF (FLDCRF-s) allows multiple latent dynamics of the class labels to interact with each other, leading to deeper understanding. FLDCRF-s outperforms CRF, LDCRF, LSTM and LSTM-CRF models in the single-label sequence tagging experiments.
- Multi-label FLDCRF (FLDCRF-m) outperforms Factorial CRF, coupled CRF, multi-label LSTM models and all considered single label sequence models (CRF, LDCRF, LSTM and LSTM-CRF) on the joint sequence tagging task. We show that while factorial CRF and coupled CRF do not yield much improvement over a CRF, the interaction among the latent dynamics in FLDCRF-m results in substantial improvements (upto 3%) over the individual LDCRF models.
- In addition to model performance on test data, we compare FLDCRF and LSTM models on several important modeling aspects, such as, ease of model selection on the validation data, consistency across validation and test performance, and computation times for training and inference (see Section 6).

The rest of the paper is organized as follows. In Section 2, we review existing literature on single and multi-label sequence modeling. In Section 3, we describe the proposed Factored Latent-Dynamic Conditional Random Fields (FLDCRF) and explain its training and inference mechanisms. In Section 4, we describe the datasets considered in our experiments. We discuss our experimental setup (feature preprocessing, models tested, metrics etc.) in Section 5 and present our results in Section 6. We briefly summarize several modeling aspects of FLDCRF and LSTM in Section 7. Finally, we offer concluding remarks and ideas for future research in Section 8.

2. Literature Review

We discuss existing literature on single and multi-label sequence modeling in this Section.

Approaches to sequence labeling can be broadly classified into two categories: a) Generative (e.g., Hidden Markov Models (Rabiner, 1989)), which learn the joint distribution $P(\mathbf{y}, \mathbf{x})$ from the training data; and b) Discriminative (e.g., Conditional Random Fields (Lafferty et al., 2001), Maximum Entropy Markov Models (McCallum et al., 2000) etc.), which directly learn the conditional distribution $P(\mathbf{y} | \mathbf{x})$ from the training data. Discriminative models have often been shown to outperform generative models on sequence labeling tasks (Lafferty et al., 2001; Morency et al., 2007; McCallum et al., 2000). In particular,

Conditional Random Fields (CRF) are frequently applied to the sequence labeling problem (Huang et al., 2015; Collobert et al., 2011; Lample et al., 2016).

A simple CRF model, also called Linear Chain CRF (LCCRF, see Fig. 1a), is defined directly over the input ($\{x_t\}_{t=1:T}$) and output ($\{y_t\}_{t=1:T}$) sequence variables. The class labels within a sequence ($\{y_t\}_{t=1:T}$) have an extrinsic dynamics (temporal dependency), as well as there exists underlying intrinsic dynamics within each class. For example, in a continuous activity recognition problem with classes ‘sitting’, ‘standing’, ‘walking’ etc., each class contains temporal transitions between certain posture states. Since such posture states are hidden, the intrinsic dynamics within a particular class can be captured by associating hidden (latent) states and allowing their transitions. The simple CRF model captures the extrinsic dynamics among the class labels through its transition links (connecting y_{t-1} and y_t). However, it does not have any latent variables in its structure and thus is unable to capture the intrinsic dynamics within the class labels.

To address this issue, Hidden-state Conditional Random Fields (HCRF, Wang et al. (2007)) introduce hidden variables in a labeled CRF structure (see Fig. 1b), in order to capture the intrinsic dynamics within the class labels. The hidden variables ($\{h_t\}_{t=t_1:t_2}$) can assume values from a predefined set of hidden states (usually specified by numbers) and the Markov dependency among the variables allows to capture the intrinsic dynamics within the class label $y_{t_1:t_2}$. However, such a structure requires prior segmentation of the training sequences according to the class labels and models each class label separately (similar to unsupervised HMMs for each class label, only trained by a discriminative approach). Hence, HCRF fails to capture the extrinsic dynamics between the class labels.

An alternative approach, named Latent-Dynamic Conditional Random Fields (LDCRF), was proposed by Morency et al. (2007) (see Fig. 1c). In a LDCRF, prior segmentation of the training sequences is not necessary. Thus, LDCRF is able to capture both intrinsic and extrinsic dynamics of the class labels. LDCRF restricts the states associated to each class label to be disjoint. Each hidden state variable h_t is also constrained to belong to the states associated to the class label y_t . These two constraints help to keep the model computations (during training and inference) tractable, despite the insertion of hidden variables. LDCRF has been shown to outperform HMM, CRF and HCRF on several sequence labeling tasks (Morency et al., 2007; Sun et al., 2008). Thai et al. (2018) recently proposed a very similar model called Embedded-State Latent CRFs. They claim to factorize the log potential as the novelty over a LDCRF, however such factorization is not reflected in their model structure and mathematical descriptions.

We propose FLDCRF-s in order to model multiple interacting latent dynamics within the class labels (see Section 1). We describe how FLDCRF-s mathematically generalizes a LDCRF in Section 3. We demonstrate improvement in performance by FLDCRF-s over LDCRF across nested cross-validation experiments on two different datasets (see Section 6).

LSTMs (Hochreiter and Schmidhuber, 1997) are the most popular kind of Recurrent Neural Networks (RNN) for sequence modeling. They are frequently applied to sequence labeling tasks, e.g., named entity recognition (Lample et al., 2016; Huang et al., 2015), noun phrase chunking (Huang et al., 2015), action recognition (Ullah et al., 2017) etc. LSTM-CRF, a variant of LSTM, combined with a CRF layer (Lample et al., 2016) for classification,

has recently been very popular for sequence labeling (Lample et al., 2016; Huang et al., 2015; Zhu et al., 2019).

A few studies (Wollmer et al., 2013; Huang et al., 2018; Zadeh et al., 2018) have compared LDCRF and LSTM on the same classification tasks. While Huang et al. (2018); Zadeh et al. (2018) show LSTM and its variants to outperform LDCRF, Wollmer et al. (2013) present results with each outperforming the other on different tasks. We show that although LSTM and LSTM-CRF outperform LDCRF on certain experiments, FLDCRF-s outperforms both LSTM and LSTM-CRF across all our experiments. By contrast, the best LSTM (and LSTM-CRF) models outshine the best FLDCRF-s performance on most of the validation sets across experiments in the paper. However, the same LSTM models (with optimized hyperparameters), kept unchanged or retrained on the validation+train data, fail to beat the optimized FLDCRF-s models on the test data across experiments (see Section 6). In addition to this inconsistent validation and test performance, LSTM models exhibit large variations in validation performance across different hyperparameter settings, with quite a few models producing significantly lesser validation performance. This variation exists without any easily discernible patterns among different hyperparameter settings. The variable performance across validation sets, lower worst cases, and inconsistency in validation and test performance make it quite tedious to select appropriate LSTM models, and raise concerns about the stability of the models in practical deployment. On the other hand, FLDCRF models with its intuitive structure, brings consistency in validation and test performance and offers a much easier model selection process. LSTM models also take significantly longer to train than FLDCRF-s models. We present results to support all the above statements in Section 6.

A factored variant of LDCRF, called multi-view LDCRF (MVLDCRF, see Fig. 6, Song et al. (2012)), was proposed to capture interactions among different latent dynamics arising due to different kinds of input features within same class labels. A MVLDCRF is a special case of FLDCRF-s, where different latent dynamics within class labels $\{y_t\}_{t=1:T}$ are learned from disjoint feature subsets $\{x_{1,t}\}_{t=1:T}$ and $\{x_{2,t}\}_{t=1:T}$ of $x_t = \{x_{1,t}, x_{2,t}\}$. Song et al. (2012) showed a MVLDCRF to outperform a LDCRF trained with the entire feature set x_t in their experiments. However, such explicit distribution of features in order to force the model to learn different latent dynamics does not guarantee improved performance, as we demonstrate by our experiments (see Section 6.3). In some cases, performance of such a MVLDCRF falls in between those of the individual LDCRFs, i.e., one trained with $x_{1,t}$ and another with $x_{2,t}$. These results are not reported by Song et al. (2012). In most other cases, a LDCRF trained with the entire feature set x_t outperforms MVLDCRF and individual LDCRFs. By contrast, a FLDCRF-s with all hidden layers modeling their dynamics from the entire x_t lets the model to learn different latent dynamics that may exist, and in general improves the LDCRF (with x_t) performance. In the case a feature distribution ($x_{1,t}$, $x_{2,t}$ etc.) is available, it is advisable to obtain individual FLDCRF-s performances on the individual feature subsets, as well as on entire x_t , and select the best performance. It is better to avoid learning from distributed features on the same training labels as in MVLDCRF, as we demonstrate on 10 different test sets in Section 6.3.

Multi-task learning (Caruana, 1997) is a very common problem of machine learning, where multiple label categories are trained and/or inferred jointly over certain input features. Modeling multiple labels together serves additional context information to individual

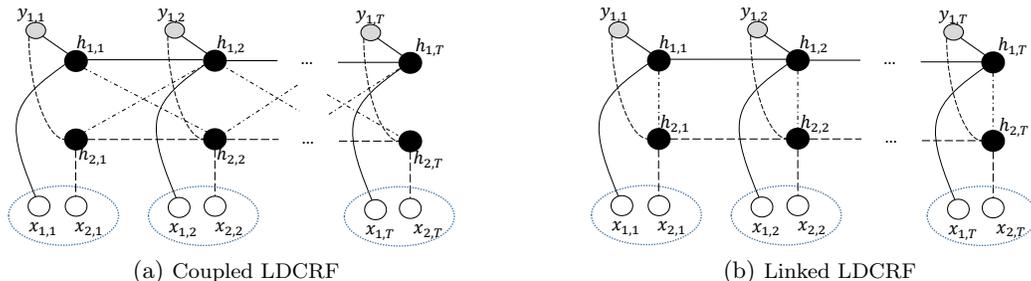


Figure 6: Multi-view LDCRFs (Song et al., 2012). The features x_t in Fig. 3b are distributed across two variables $x_{1,t}$ and $x_{2,t}$, i.e., $x_t = \{x_{1,t}, x_{2,t}\}$.

labeling tasks, which help to improve the overall labeling performance. However, the model must allow appropriate label interactions for beneficial results. Multi-task learning is frequently applied to tag images with multiple labels (Wang et al., 2016; Wei et al., 2015).

Multi-task sequence learning (Sutton et al., 2007; Collobert et al., 2011) is the task to jointly tag sequence values with multiple label categories. Existing multi-task sequence labeling approaches can be broadly classified into two categories:

1. **Using a joint state variable** (see Fig. 4a), i.e., $y_t = \{y_{1,t}, y_{2,t}\}$ (Li et al., 2015; Luo et al., 2015; Pahuja et al., 2017). Li et al. (2015) learned a coupled CRF (CCRF) model that allowed the two (or more) label categories to interact temporally through their coupled transition matrix. Luo et al. (2015) proposed a model called JERL (Joint Entity Recognition and Linking) for joint named entity recognition and disambiguation, which is similar to a coupled CRF. A correlated BiRNN model was introduced by Pahuja et al. (2017), where multiple label categories (punctuation and capitalization) were combined to pass through the same RNN hidden layer, thus capturing the multi-label interactions through the temporal transitions of the hidden layer. Although all such joint models showed improvements over the individual models, the joint state space makes the model too large with multiple classes (or latent states) for each task.
2. **Using separate state variables** ($y_{1,t}$, $y_{2,t}$ etc.) **with cotemporal links** (see Fig. 4b-c Sutton et al. (2007)) to capture interaction. These structures are derived from Dynamic CRFs proposed by Sutton et al. (2007). The structure in Fig. 4b is also well known as a Factorial CRF (FCRF). These factored representations from a DCRF require fewer parameters than modeling with the joint state variable, and allow easy addition/removal of links depending on required model complexity and available training data. For instance, the coupled CRF model by (Li et al., 2015) can be efficiently represented by Fig. 4c, with substantially reduced model parameters for the same training labels $\{y_{1,t}\}_{t=1:T}$ and $\{y_{2,t}\}_{t=1:T}$. Sutton et al. (2007) showed improvements in labeling accuracy with a FCRF while considering a joint noun phrase (NP) chunking and part of speech (POS) tagging task. They also reported better performance for FCRF compared to two separate CRFs for reduced training data.

Several other studies (Collobert et al., 2011; Changpinyo et al., 2018; Dredze et al., 2009; Shi et al., 2007; Bruce et al., 2015) have been conducted on joint sequence labeling. Shi et

al. (2007) proposed a dual layer CRF for joint decoding of different tasks (segmentation and POS tagging), however showed marginal improvement over individual labeling with their approach. Dredze et al. (2009) designed a Multi-CRF with a shared entropy based training likelihood function. Collobert et al. (2011) and Changpinyo et al. (2018) reviewed existing approaches to multi-task sequence learning and argued that joint learning and decoding may not necessarily give better results than individual approaches. Collobert et al. (2011) also questioned the availability of fully labeled datasets as a bottleneck of such approaches. As discussed earlier, models represented by distributed state variables (see Fig. 4b-c) are quite flexible and can be trained on sequences with missing labels by simply adding/dropping links whenever necessary, depending on the variable (label node) availability. It is hard to tackle such data for models represented by the joint state variable (see Fig. 4a), as it requires both label categories $y_{1,t}$, $y_{2,t}$ to form the joint variable y_t for training at all t . While most multi-task sequence learning models have been designed for problems in Natural Language Processing (NLP), Bruce et al. (2015) presented a joint model to improve the performance of action recognition and pose estimation from video data.

As elaborated in Section 1, we propose FLDCRF-m in order to model latent-dynamic interactions among different label categories. We describe how FLDCRF-m gives a generic expression for DCRFs with embedded hidden-state variables in Section 3. FLDCRF-m outperforms all state-of-the-art single (CRF, LDCRF, LSTM, LSTM-CRF) and multi-label (FCRF, CCRF, LSTM-m) models in the multi-label sequence tagging experiment.

Most studies on single and joint sequence labeling have been conducted in the context of NLP, viz., CoNLL 2000 chunking (Sang and Buchholz, 2000), CoNLL 2003 named entity tagging (Sang and Meulder, 2003) etc. However, recent reported (F1-measure) performance results on these tasks are largely driven by the quality of the input features, while the models do not show much variance in performance. By contrast, we wish to determine the effectiveness of the proposed models and therefore do not consider such datasets for evaluating our models. We plan to apply FLDCRF on NLP tasks in future.

We apply the proposed FLDCRF-s and FLDCRF-m models to two problems:

1. Continuous gesture recognition on the UCI gesture phase segmentation dataset (Madedo et al., 2013). The task concerns online segmentation of gestures from rest positions.
2. Continuous multi-action recognition on the UCI opportunity dataset (Chavarriaga et al., 2013). This task aims at continuous joint recognition of two action types (locomotion and a high-level activity, see Section 4).

We refer to Sections 4 and 5 for detailed description of the datasets and experiments respectively. All experiments in this paper consider continuous online recognition problems, i.e., the inference tasks are formulated in terms of $P(y_t | x_{1:t})$. We do not utilize future input features (x_{t+1} , x_{t+2} etc.) and future label/state variables during modeling/testing in any of the models.

3. Factored Latent-Dynamic Conditional Random Fields (FLDCRF)

We describe the proposed FLDCRF model here. As described in Section 1, a fully labeled FLDCRF has three variants:

1. FLDCRF single label sequence model (FLDCRF-s, see Fig. 3b),
2. FLDCRF multi label sequence model (FLDCRF-m, see Fig. 5a), and
3. FLDCRF multi agent interaction model (FLDCRF-i, see Fig. 5b).

All 3 three models can be expressed via a similar mathematical representation. Since the FLDCRF-s and FLDCRF-i are special cases of the FLDCRF-m model, we describe the FLDCRF-m model for better generalization. We describe how we construct the FLDCRF model, give its generic expression and then introduce the simplified versions we apply in this paper.

3.1 Model

Fig. 5a shows the graph structure for FLDCRF in multi-label classification problems (FLDCRF-m). We depict the model for $L = 2$ different label categories $\{y_{1,t}\}_{t=1:T}$ and $\{y_{2,t}\}_{t=1:T}$, interacting with each other through their respective hidden layers $\{h_{1,t}\}_{t=1:T}$ and $\{h_{2,t}\}_{t=1:T}$. In the case of a FLDCRF single-label model (FLDCRF-s, see Fig. 3b), the label category $\{y_t\}_{t=1:T}$ is associated to all hidden layers.

FLDCRF-m model (see Fig. 5a) can be easily extended to accommodate multiple hidden layers for each label category $\{y_{i,t}\}_{t=1:T}$, $i = 1 : L$. However, for simplicity of model description, we assume only one hidden layer per label category in FLDCRF-m. This leads to a total L hidden layers in the model described below. We also assume first-order Markov connections among different hidden layers while applying FLDCRF in this paper. Higher order connections are possible but avoided for simplicity. These extensions can be utilized in order to improve model performance at the cost of model complexity and more training data, however we leave them to user discretion.

Let, $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$ denote the sequence of observations. $\mathbf{y}_i = \{y_{i,1}, y_{i,2}, \dots, y_{i,T}\}$ are the observed labels along layer i , $i = 1 : L$. In the case of single-label prediction task by a FLDCRF-s (see Fig. 3b), all layers take the same labels during model training, i.e., $y_{i,t}$ is same $\forall i$. In the case of multi-agent interaction model (FLDCRF-i, see Fig. 5b), the features x_t are distributed across L layers, i.e., $x_t = \{x_{i,t}\}_{i=1:L}$, assuming L interacting agents in the environment.

Let, Υ_i be the alphabet for all possible label categories in layer i , $i = 1 : L$, i.e., $y_{i,t} \in \Upsilon_i$. $\mathbf{h}_i = \{h_{i,1}, h_{i,2}, \dots, h_{i,T}\}$ constitutes the i -th hidden layer. Each possible class label $\ell_i \in \Upsilon_i$ in layer i is associated with a set of hidden states \mathcal{H}_{i,ℓ_i} . \mathcal{H}_i is the set of all possible hidden states for layer i given by $\mathcal{H}_i = \bigcup_{\ell_i} \mathcal{H}_{i,\ell_i}$.

The joint conditional model is defined as:

$$P(\{\mathbf{y}_i\}_{1:L} | \mathbf{x}, \theta) = \sum_{\{\mathbf{h}_i\}_{1:L}} P(\{\mathbf{y}_i\}_{1:L} | \{\mathbf{h}_i\}_{1:L}, \mathbf{x}, \theta) \cdot P(\{\mathbf{h}_i\}_{1:L} | \mathbf{x}, \theta). \quad (1)$$

In order to keep model computations during training and testing tractable, we introduce the layers of hidden variables $\{\mathbf{h}_i\}_{1:L} = \{h_{i,t}\}_{i=1:L,t=1:T}$ to the model with links (graphical

constraints) as depicted in Fig. 7. This allows us to factorize $P(\{\mathbf{y}_i\}_{1:L} \mid \{\mathbf{h}_i\}_{1:L}, \mathbf{x}, \theta)$ according to equation (2).

$$P(\{\mathbf{y}_i\}_{1:L} \mid \{\mathbf{h}_i\}_{1:L}, \mathbf{x}, \theta) = \prod_{i=1}^L \prod_{t=1}^T P(y_{i,t} \mid h_{i,t}). \quad (2)$$

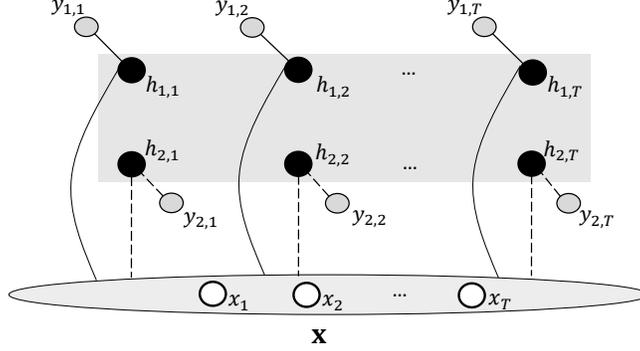


Figure 7: The general structure of FLDCRF with the essential graphical constraints. Hidden variables $h_{i,t}$, $i = 1 : L$, $t = 1 : T$ can assume any dependency structure (first-order Markov, second-order Markov, cotemporal links between layers etc.) within themselves in this general FLDCRF structure.

Replacing $P(\{\mathbf{y}_i\}_{1:L} \mid \{\mathbf{h}_i\}_{1:L}, \mathbf{x}, \theta)$ from equation (2), we can re-write equation (1) as:

$$\begin{aligned} P(\{\mathbf{y}_i\}_{1:L} \mid \mathbf{x}, \theta) &= \sum_{\{\mathbf{h}_i\}_{1:L} : \forall h_{i,t} \in \mathcal{H}_{i,y_{i,t}}} \left(\prod_{i=1}^L \prod_{t=1}^T P(y_{i,t} \mid h_{i,t}) \right) \cdot P(\{\mathbf{h}_i\}_{1:L} \mid \mathbf{x}, \theta) \\ &+ \sum_{\{\mathbf{h}_i\}_{1:L} : \exists h_{i,t} \notin \mathcal{H}_{i,y_{i,t}}} \left(\prod_{i=1}^L \prod_{t=1}^T P(y_{i,t} \mid h_{i,t}) \right) \cdot P(\{\mathbf{h}_i\}_{1:L} \mid \mathbf{x}, \theta). \end{aligned} \quad (3)$$

In order to further reduce model computations, we define the following model constraints, extending on Morency et al. (2007):

1. \mathcal{H}_{i,ℓ_i} are disjoint $\forall \ell_i \in \Upsilon_i$, $\forall i = 1 : L$.
2. $h_{i,t}$ can only assume values from the set of hidden states assigned to the label $y_{i,t}$, i.e., $h_{i,t} \in \mathcal{H}_{i,y_{i,t}}$, $\forall i = 1 : L$ and $\forall t = 1 : T$.

These constraints let us write the following:

$$P(y_{i,t} = \ell_i \mid h_{i,t}) = \begin{cases} 1, & h_{i,t} \in \mathcal{H}_{i,y_{i,t}=\ell_i} \\ 0, & h_{i,t} \notin \mathcal{H}_{i,y_{i,t}=\ell_i}. \end{cases} \quad (4)$$

Thus, the FLDCRF model in equation (3) can be reduced to:

$$P(\{\mathbf{y}_i\}_{1:L} \mid \mathbf{x}, \theta) = \sum_{\{\mathbf{h}_i\}_{1:L} : \forall h_{i,t} \in \mathcal{H}_{i,y_{i,t}}} P(\{\mathbf{h}_i\}_{1:L} \mid \mathbf{x}, \theta). \quad (5)$$

Equation (5) gives the general expression of a FLDCRF.

Equation (5) simplifies to a LDCRF (Morency et al., 2007) model for $L = 1$. If we assume only one distinct hidden state per class label in each layer i , i.e., $|\mathcal{H}_{i,l_i}| = 1, \forall i, \forall l_i$; and $|\mathcal{H}_i| = |\Upsilon_i|, \forall i$; equation (5) simplifies to the joint conditional distribution $P(\{\mathbf{y}_i\}_{1:L} | \mathbf{x}, \theta)$ given by a DCRF (Sutton et al., 2007) or a coupled CRF (Li et al., 2015). It is also straightforward to see that FLDCRF yields a LCCRF (Lafferty et al., 2001) when $L = 1$ and one distinct hidden state is associated to each class label. Thus, FLDCRF subsumes major sequential CRF variants viz., LDCRF, DCRF and LCCRF (illustrated in Fig. 8).

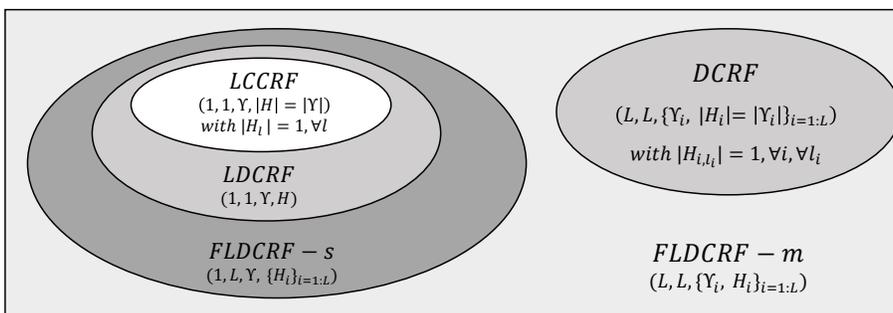


Figure 8: Venn diagram illustrating relationship between FLDCRF and major sequential CRF variants. The configuration underneath the model names describe (number of label categories, number of hidden layers, class alphabet and the sets of hidden states) in a FLDCRF-m model in order to derive it. For simplicity of description, we assume 1 hidden layer for each label category in FLDCRF-m configuration, i.e., L hidden layers for L label categories.

We define $P(\{\mathbf{h}_i\}_{1:L} | \mathbf{x}, \theta)$ by the standard CRF formulation,

$$P(\{\mathbf{h}_i\}_{1:L} | \mathbf{x}, \theta) = \frac{1}{\mathbf{Z}(\mathbf{x}, \theta)} \exp\left(\sum_k \theta_k \cdot F_k(\{\mathbf{h}_i\}_{1:L}, \mathbf{x})\right), \quad (6)$$

where index k ranges over all parameters $\theta = \{\theta_k\}$ and $\mathbf{Z}(\mathbf{x}, \theta)$ is the partition function given by:

$$\mathbf{Z}(\mathbf{x}, \theta) = \sum_{\{\mathbf{h}_i\}_{1:L}} \exp\left(\sum_k \theta_k \cdot F_k(\{\mathbf{h}_i\}_{1:L}, \mathbf{x})\right). \quad (7)$$

Now we introduce first-order Markov assumptions (as described earlier) among the hidden variables $\{\mathbf{h}_i\}_{1:L}$. Therefore, the feature functions F_k 's are factorized (i.e., summed inside exponent) as:

$$F_k(\{\mathbf{h}_i\}_{1:L}, \mathbf{x}) = \sum_{t=1}^T f_k(\{h_{i,t-1}\}_{1:L}, \{h_{i,t}\}_{1:L}, \mathbf{x}, t). \quad (8)$$

We further factorize each $f_k(\{h_{i,t-1}\}_{1:L}, \{h_{i,t}\}_{1:L}, \mathbf{x}, t)$ at slice t , such that each hidden layer assumes its own (temporal) dynamics over the observations, while interacting with each other through cotemporal and first-order Markov (omitted in figures to avoid clutter) *influence* links. Thus, each factorized component of $f_k(\{h_{i,t-1}\}_{1:L}, \{h_{i,t}\}_{1:L}, \mathbf{x}, t)$ at slice t can be a *state* function $s_k(h_{i,t}, \mathbf{x}, t)$, a *transition* function $t_k(h_{i,t-1}, h_{i,t}, \mathbf{x}, t)$ or an *influence* function $i_k(h_{i,t-l}, h_{j,t-m}, \mathbf{x}, t)$, with $i, j \in \{1 : L\}$ and $(l, m) \in \{(0, 0), (0, 1), (1, 0)\}$. We define *state* and *transition* functions by the following indicator functions:

$$\begin{aligned} s_k(h_{i,t}, \mathbf{x}, t) &= \mathbb{1}_{\{(h_{i,t})=k\}} \cdot x_t, \\ t_k(h_{i,t-1}, h_{i,t}, \mathbf{x}, t) &= \mathbb{1}_{\{(h_{i,t}, h_{i,t-1})=k\}}, \end{aligned} \tag{9}$$

assuming one-hot representations for discrete x_t components. We define *influence* functions $i_k(h_{i,t-l}, h_{j,t-m}, \mathbf{x}, t)$, with $i, j \in \{1 : L\}$ and $(l, m) \in \{(0, 0), (0, 1), (1, 0)\}$ as:

$$i_k(h_{i,t-l}, h_{j,t-m}, \mathbf{x}, t) = \mathbb{1}_{\{(h_{i,t-l}, h_{j,t-m})=k\}}. \tag{10}$$

For the interaction model depicted in Fig. 5b, the mathematical expressions are identical, only with a minor change in the *state* function:

$$s_k(h_{i,t}, \mathbf{x}_i, t) = \mathbb{1}_{\{(h_{i,t})=k\}} \cdot x_{i,t}, \tag{11}$$

where $\mathbf{x}_i = \{x_{i,t}\}_{t=1:T}$.

3.2 Training

We estimate the model parameters by maximizing the conditional log-likelihood of the training data:

$$\mathbf{L}(\theta) = \sum_{n=1}^N \log P(\{\mathbf{y}_i\}_{1:L}^{(n)} \mid \mathbf{x}^{(n)}, \theta) - \frac{\|\theta\|^2}{2\sigma^2}, \tag{12}$$

where N is the total number of available labeled sequences. A L2 regularizer (second term in equation (12)) was included in order to reduce overfitting during our experiments.

$P(\{\mathbf{y}_i\}_{1:L}^{(n)} \mid \mathbf{x}^{(n)}, \theta)$, $n = 1 : N$ are obtained from equations (5)-(10). The numerator $\mathbf{N}(\mathbf{x}, \theta) = \sum_{\{\mathbf{h}_i\}_{1:L} : \forall h_{i,t} \in \mathcal{H}_{i,y_{i,t}}} \exp\left(\sum_k \theta_k \cdot F_k(\{\mathbf{h}_i\}_{1:L}, \mathbf{x})\right)$ and the denominator $\mathbf{Z}(\mathbf{x}, \theta)$ of equation (5), given by (6) and (7), are in the classic sum-product form of dynamic programming and are efficiently computed by the forward algorithm (Rabiner, 1989). We apply the default BFGS optimizer in Stan modeling language (Carpenter et al., 2017) to obtain the estimates $\hat{\theta} = \{\hat{\theta}_k\}$.

3.3 Inference

Multiple label sequences \mathbf{y}_i , $i = 1 : L$, can be inferred from the same graph structure by marginalizing over other labels:

$$\hat{\mathbf{y}}_i = \operatorname{argmax}_{\mathbf{y}_i} \sum_{\{\mathbf{y}_i\}_{1:L-\mathbf{y}_i}} P\left(\{\mathbf{y}_i\}_{1:L} \mid \mathbf{x}, \hat{\theta}\right), \quad (13)$$

where $\mathbf{x} = \{x_t\}_{t=1:T}$ is the observed input sequence of length T . $P\left(\{\mathbf{y}_i\}_{1:L} \mid \mathbf{x}, \hat{\theta}\right)$ can be obtained from (5) and estimated parameters $\hat{\theta}$. At each instant t , the marginals $P(\{h_{i,t}\}_{1:L} \mid \mathbf{x}, \hat{\theta})$ are computed and summed according to the disjoint sets of hidden states to obtain joint estimates of desired labels $\hat{y}_{i,t}$, $t = 1, 2, \dots, \forall i = 1 : L$, as follows:

$$P(\{y_{i,t}\}_{1:L} \mid \mathbf{x}) = \sum_{\{h_{i,t}\}_{1:L}: h_{i,t} \in \mathcal{H}_{i,y_{i,t}}} P\left(\{h_{i,t}\}_{1:L} \mid \mathbf{x}, \hat{\theta}\right). \quad (14)$$

After marginalizing according to (13), the label $\hat{y}_{i,t}$ corresponding to the maximum probability is inferred. Since we consider online inference (i.e., input sequence observed upto current instant t) for our continuous gesture and action recognition problems, the inference problem in equation (14) gets reduced to $P(\{y_{i,t}\}_{1:L} \mid x_{1:t})$. We compute the necessary probabilities $P\left(\{h_{i,t}\}_{1:L} \mid x_{1:t}, \hat{\theta}\right)$ by the forward algorithm (Rabiner, 1989). Forward-backward algorithm (Rabiner, 1989) and Viterbi algorithm (Forney, 1973) can also be applied for problems where online inference is not necessary.

4. Datasets

We apply our models on two datasets: a) UCI gesture phase segmentation dataset (Madeo et al., 2013), and b) UCI opportunity dataset (Chavarriaga et al., 2013). We describe the datasets below.

4.1 UCI Gesture Phase Dataset

The UCI gesture phase data consists of features extracted from 7 videos with people gesticulating while telling stories. The data is captured by a Microsoft Xbox Kinect. Each time point within the 7 sequences is described by 18 positional features, 32 dynamic features (velocity, acceleration etc.) and is labeled with one of the five gesture classes: rest, preparation, hold, stroke and retraction. The number of instances in each sequence is tabulated in Table 1. ‘A’, ‘B’, and ‘C’ denote different participants and ‘1’, ‘2’, ‘3’ describe different stories. Although the data is multiclass, the most popular studies on the dataset (Madeo et al., 2013; Wagner et al., 2014) have only considered binary classification, i.e., segmentation of gesture positions (i.e., positives - preparation, hold, stroke, retraction) from the rest (negative) position. We perform two binary time-series continuous gesture segmentation experiments on this dataset:

1. We follow the experiment 2 in (Wagner et al., 2014). The goal of the experiment is to train on the sequence A1 and test on sequence A2. We utilize the ‘Data vector 2’ in the paper as our input features x_t , which consists of 3 dimensional positions of the left hand, right hand, left wrist and right wrist together, yielding a 12 dimensional input feature set. During experiment, we selected our models (hyperparameters, see

Section 5) on a validation set comprising of the last 30% frames of sequence A1, while training on the first 70%. During testing, we re-trained the selected models on the entire A1 sequence.

2. A nested cross validation, with 7 outer loops (one for testing each sequence) and 6 inner loops (leave-one-out) per outer loop. We present results of both our experiments on the UCI gesture phase data in Section 6.1.

Table 1: Table showing number of instances in each sequence in UCI gesture phase dataset (Madeo et al., 2013).

Sequence	A1	A2	A3	B1	B3	C1	C3	Total
Num frames	1747	1264	1834	1073	1423	1111	1448	9900

4.2 UCI Opportunity Dataset

The UCI opportunity dataset (Chavarriaga et al., 2013) is a public activity recognition dataset. The data contains 20 ADL (activity of daily living) sequences from 4 different participants (S1, S2, S3 and S4). Each participant has 5 ADL sequences. Each instance in the data has 3 activity label categories:

- locomotion, which has four classes, viz., ‘lie’, ‘sit’, ‘stand’ and ‘walk’;
- a high-level activity (HL), which has five classes, viz., ‘relaxing’, ‘early morning’, ‘coffee time’, ‘sandwich time’ and ‘clean up’; and
- a mid-level activity (ML), which has several classes.

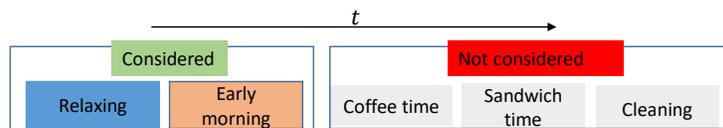


Figure 9: Illustrating each ADL sequence in the UCI opportunity dataset (Chavarriaga et al., 2013).

We do not utilize the mid-level activity labels in this paper. The 5 high-level activities are performed in sequence within each complete ADL run (see Fig. 9). We consider the 5 sequences (ADL runs) from participant S2 in our experiment. From each ADL sequence, we utilize the data upto two high-level activity classes (i.e., ‘Relaxing’ and ‘Early morning’), illustrated in Fig. 9. The ‘Relaxing’ portions of the sequence primarily consists of the ‘lie’ and ‘sit’ locomotion classes, while the ‘early morning’ activity class comes mostly with ‘stand’ and ‘walk’ locomotion classes, making the two label categories depending on each other. Each time-point in the data also contains 242 sensor outputs (113 body-worn sensors, 32 shoe sensors and 97 other object and ambient sensors). We utilize the body-worn and shoe sensor outputs as our input features, yielding a 145 dimensional input feature set. For

our convenience, we further divided the data into 28 sub-sequences, each fully labeled by the two label categories. Additionally, there are instances within the data with missing values for some sensor outputs. We replaced such values with the previous instant, if available, otherwise by ‘0’s. Since the FLDCRF models encode the input features (x_t) by exponential functions over ($\theta_k \cdot x_t$), replacing values by ‘0’s means a multiplication by 1, thus not affecting the model likelihood. We preprocessed each input dimension of x_t between 0 to 1 for faster learning of our models.

We perform a nested cross-validation with 5 outer loops (with test sequences coming from each ADL sequence) and 4 inner loops (leave-one-ADL sequence out) per outer loop. Tabulated data (features and labels) for our multi-label sequence tagging experiment on UCI opportunity data is available here: <https://github.com/satyajitneogiju/FLDCRF-for-sequence-labeling>. Table 2 provides the details (length, ADL run etc.) of the 28 sub-sequences considered in this experiment.

Important dataset characteristics of the two datasets are listed in Table 3.

Table 2: Number of instances in each sequence in the nested CV experiment on UCI opportunity dataset (Chavarriaga et al., 2013).

Nested outer set \ Sequence	1	2	3	4	5	6	7	8	9	Total
1	1000	1000	800	943	724	847	1352	1276	298	8240
2	721	863	800	1428	216					4028
3	583	843	800	746	655					3627
4	477	253	1008	954	583					3275
5	711	1338	626	493						3168
Total										22338

5. Experimental Setup

5.1 Metrics for evaluation

We assess the different models by the F1 score. For the gesture phase segmentation problem on UCI gesture phase dataset, we consider rest position as ‘negative’ and other gestures as ‘positive’. The F1 measure (F) is computed as:

$$F = \frac{2 \cdot P \cdot R}{P + R}, \quad (15)$$

where P and R represent precision and recall, defined as $P = \frac{TP}{TP+FP}$ and $R = \frac{TP}{TP+FN}$ respectively. TP , FP and FN are the predicted true positives, false positives and false negatives by the model.

Table 3: Information on the UCI gesture phase (Madeo et al., 2013) and the UCI opportunity dataset (Chavarriaga et al., 2013).

Attribute \ Dataset	UCI gesture phase	UCI opportunity
Type	Single label	Multi label
No. sequences	7	5
No. sub-sequences	7	28
Experiments	Experiment 2 (Wagner et al., 2014) Nested LOOCV	Nested LOOCV
No. labels	1 (Gesture)	2 (Locomotion & High-level activity)
No. classes/label	2	4 & 2
Tested models	CRF, LDCRF, FLDCRF-s, LSTM, LSTM-CRF	CRF, LDCRF, FLDCRF-s, LSTM, LSTM-CRF, FCRF, CCRF, FLDCRF-m, LSTM-m
Input dimension ($ x_t $)	12	145
Metrics	F1-score	Micro F1-score

For the multi-class action recognition problem on UCI opportunity data, we assess the models by the micro F1 score (F_{micro}), which also is defined by (15), where precision and recall are given by:

$$P_{micro} = \frac{\sum_{c=1}^{N_c} TP_i}{\sum_{c=1}^{N_c} TP_i + \sum_{c=1}^{N_c} FP_i}, \quad R_{micro} = \frac{\sum_{c=1}^{N_c} TP_i}{\sum_{c=1}^{N_c} TP_i + \sum_{c=1}^{N_c} FN_i}. \quad (16)$$

We apply the micro F1-score (F_{micro}) for evaluating the ‘locomotion’ (4 classes) and overall (6 classes) performance. As both class labels (‘relaxing’ and ‘early morning’) of the high-level activity (HL) are equally important, we slightly modify the precision (P_{HL}) and recall (R_{HL}) as follows: $P_{HL} = \frac{TP_{early} + TP_{relax}}{TP_{early} + TP_{relax} + FP_{early}}$ and $R_{HL} = \frac{TP_{early} + TP_{relax}}{TP_{early} + TP_{relax} + FN_{early}}$.

During validation, we select the multi-label sequential models (FLDCRF-m and LSTM-m) based on their combined performance (on all 6 classes) in the inner loops. The single-label models are separately selected on each label category.

5.2 Benchmarking

For the single label sequence tagging tasks (‘experiment 2’ and nested CV) on the UCI gesture phase data, we compare CRF, LDCRF, FLDCRF-s, LSTM, and LSTM-CRF. For the multi-label sequence tagging task, we compare CRF, LDCRF, FLDCRF-s, LSTM, LSTM-CRF, FCRF, CCRF, FLDCRF-m, and LSTM-m.

A FLDCRF-s model has two hyperparameters: number of hidden layers (N_h) and number of hidden states per label ($\{N_{si}\}$) along layers $i = 1 : L$. For simplicity, we apply the

same number of hidden states along all layers, i.e., $N_{si} = N_s, \forall i = 1 : L$. We denote such model by FLDCRF-s($\langle N_h \rangle / \langle N_s \rangle$). For example, if we consider 2 hidden layers and 3 hidden states (per class label) across each layer, then the FLDCRF-s model will be denoted by FLDCRF-s(2/3). In Table 4, we provide the list of tested hyperparameter settings of FLDCRF-s across the 3 different experiments.

A FLDCRF-m (multi-label) model has three hyper-parameters: number of different label categories N_l , number of hidden layers per label category $\{N_{hl}\}_{l=1:N_l}$, and number of hidden states per label $\{N_{si}\}$ along layers $i = 1 : L$, where $L = \sum_{l=1}^{N_l} \{N_{hl}\}$. For simplicity, we only keep one hidden layer ($N_{hl} = 1, \forall l = 1 : N_l$) for each label category, thus making total number of hidden layers equal to the number of different label categories, i.e., $L = N_l$. We denote such a model by FLDCRF-m $\{N_{si}\}_{i=1:L}$. For example, if there are 3 label categories and we associate 1, 2 and 4 hidden states (per class label) respectively to each label category, then the 3-layered FLDCRF-m model will be represented as FLDCRF-m $\{1, 2, 4\}$. In our multi-label experiment on the opportunity data, there are two different label categories (locomotion and high-level), and the FLDCRF-m models are denoted as FLDCRF-m $\{N_{s1}, n_{s2}\}$. In Table 4, we list the considered hyperparameter settings of FLDCRF-m. We considered two different FLDCRF-m models (with and without the first-order Markov influence among hidden layers, see Fig. 10) during our multi-label experiment on the UCI opportunity data.

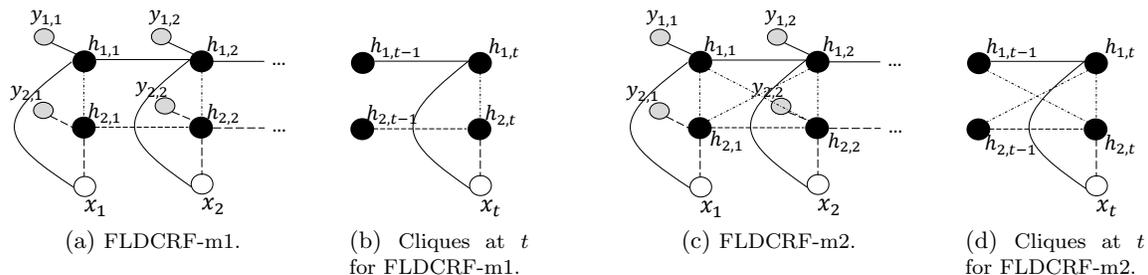


Figure 10: Two different FLDCRF-m models considered in the experiment on the UCI opportunity data.

An LSTM has three hyperparameters: number of hidden units (N_{hls}), number of epochs to train (N_{els}), and training minibatch size (N_{mls}). To the best of our knowledge, there is no generally accepted rule to select the optimal N_{hls} , therefore, they need to be tuned by trial and error. As suggested in (Eckhardt, 2018; Thomas, 2017; Stack exchange, 2018), we vary the N_{hls} according to,

$$N_{hls} = \frac{N_{sa}}{\alpha \cdot (N_i + N_o)}, \tag{17}$$

where N_{sa} is the number of samples in the training dataset (number of instances; see Tables 1 and 2), N_i is the number of input neurons ($|x_t|$ in our case; see Table 3) and N_o is the number of output neurons from the LSTM layer. We vary the parameter α between 2 and 10 to select the N_{hls} 's. A few other popular suggestions include:

- Keep N_{hls} between N_i and N_o ,

- Set N_{hls} as the arithmetic mean of N_i and N_o ,
- Set N_{hls} as the geometric mean of N_i and N_o etc.

We try to follow these recommendations as closely as possible while selecting N_{hls} values to be tested. In Table 4, we list the considered LSTM N_{hls} values in each of our experiments. During validation, we run each model setting (N_{hls}) for 500 epochs, saving model performances on validation sets at every 100 epochs. We feed all training sequences to the LSTM models at each epoch at the rate of 1 sequence per minibatch.

Table 4: Hyperparameter settings for FLDCRF-s, FLDCRF-m, LSTM and LSTM-m.

Experiment	Models			
	FLDCRF-s (N_h/N_s)	FLDCRF-m ($\{N_{s1}, n_{s2}\}$)	LSTM (N_{hls})	LSTM-m (N_{hls})
Experiment 2 (Wagner et al., 2014)	1/1, 1/2, 1/3, 1/4, 1/5, 1/6, 1/7, 2/1, 2/2, 2/3, 2/4, 3/1, 3/2, 3/3	NA	2, 5, 7, 10, 14, 20, 50, 100, 200	NA
Nested CV on Gesture phase	1/1, 1/2, 1/3, 1/4, 1/5, 1/6, 2/1, 2/2, 2/3, 2/4, 2/5, 3/1, 3/2, 3/3	NA	2, 5, 7, 10, 14, 20, 50, 100, 200	NA
Nested CV on Opportunity (Locomotion)	1/1, 1/2, 1/3, 1/4, 2/1, 2/2, 2/3	$\{1,1\}, \{1,2\}, \{1,3\},$ $\{1,4\}, \{1,5\}, \{1,6\},$ $\{2,1\}, \{2,2\}, \{2,3\},$ $\{2,4\}, \{2,5\}, \{2,6\},$ $\{3,1\}, \{3,2\}, \{3,3\}$	5, 10, 25, 50, 75, 150, 300, 500	5, 10, 25, 50, 75, 150, 300, 500
Nested CV on Opportunity (HL)	1/1, 1/2, 1/3, 1/4, 1/5, 1/6, 2/1, 2/2, 2/3, 2/4		2, 5, 10, 25, 50, 75, 150	

5.3 Model Platforms and System Specifications

We trained and tested FLDCRF type models (FLDCRF, LDCRF, CRF, CCRF etc.) by the PyStan interface of Stan modeling language (Carpenter et al., 2017). Stan’s in-built BFGS and L-BFGS optimizers allow us to conveniently train models by defining the model likelihoods. We apply BFGS to train all FLDCRF models and keep all default settings for the optimizer.

We implemented the LSTM models defined in Keras, a deep learning library for Python running Tensorflow in the backend. LSTM parameters were trained by an Adam optimizer with a default learning rate of 0.001 and other default parameters. To reduce overfitting, we added a dropout layer (with regularization 0.2) between the LSTM and softmax (or

CRF for LSTM-CRF model) layers for all models with ≥ 5 hidden units. We also added a dropout (with regularization 0.2) to the LSTM layer for all models with ≥ 5 hidden units. Our Pystan FLDCRF and Keras LSTM implementation codes are available from - <https://github.com/satyajitneogiju/FLDCRF-for-sequence-labeling>.

We perform the metric computations and plot figures in MATLAB 2015b (Mathworks Corp., 2015). The LSTM models are trained (and tested) on an Nvidia Tesla K80 GPU, while the FLDCRF models are trained (and tested) on an Intel(R) Xeon(R) CPU E5-2630 v3 @2.40GHz CPU.

6. Results

A good model should not only perform well on the validation data, it must be consistent across validation and test data. At the same time, the model selection process should be lucid and easy. It is also desirable to have fast model training and inference mechanisms, preferably without additional resources, e.g., GPU. Therefore, we not only compare different models (FLDCRF, LSTM, LSTM-CRF etc.) on the test data, but also examine several other modeling aspects, viz., ease of model selection, consistency across validation and test performance, and computation times. We present the following 4 types of model performance measures for each experiment:

- *Test performance:* Under this measure, different models are compared on the test data.
- *Ease of model selection:* Under this performance attribute, we discuss the model hyperparameter selection process of FLDCRF and LSTM families on the validation data. In order for this model selection process to be fast, easy and effective, a model must have the following characteristics:
 1. Fewer hyperparameters: There should be very few types of hyperparameters which strongly influence the model performance. FLDCRF models only depend on the N_h/N_s settings, and do not need to be tuned for number of training epochs. LSTM models, on the other hand, need to be optimized for number of training epochs alongside the number of hidden units N_{hls} and training minibatch size N_{mls} .
 2. Rule to choose hyperparameters: It is best to have definite rules to choose the model hyperparameters, rather having to select them via trial and error.
 3. Observable pattern in validation performance: In order to easily select the model hyperparameters, the model must show discernible patterns (increasing/decreasing) in validation performance against the choices of the hyperparameters. Model performance across different hyperparameter settings should not vary at random.
 4. Lesser variance: In order to avoid tedious hyperparameter selection and continuous monitoring of model performance on validation data, model performance across adjacent hyperparameter settings should not vary widely.

Wide and random variation across hyperparameters (on validation data) without any pattern also require widespread selection of hyperparameters.

5. Fair worst case performance: The model must guarantee good performance with minimal effort to optimize hyperparameters. To accomplish this, worst case performances reported by the model on validation data should not differ largely from the best/average performances, as well as percentage of such poor outcomes should be very low (if any).

- *Consistency*: We examine consistency of the selected models on test data under this measure. Consistency in model performance across validation and test data reflects stability of a model for practical deployment.
- *Computation times*: In this result category, we present computation times required by the models for training and inference, for different choices of the hyperparameters.

6.1 UCI Gesture Phase Dataset

As mentioned earlier, we perform two experiments (see Section 4.1) on the UCI Gesture Phase Dataset for segmenting gestures from rest positions. In the first experiment, we follow the experiment 2 by Wagner et al. (2014). In the second experiment, we perform a 7 (outer) fold nested cross-validation.

6.1.1 EXPERIMENT 1: (WAGNER ET AL., 2014)

- *Test performance*:

Table 5 compares different models on the test sequence A2 of the UCI gesture phase dataset. FLDCRF-s considerably outperforms LSTM and LSTM-CRF models on the test set. FLDCRF-s also outshines a multi-layered perceptron (MLP) model (Wagner et al., 2014) reported on the test data. However, FLDCRF-s does not improve LDCRF performance in this experiment.

Considering variable performance on re-training the LSTM models, we also report performance of the optimized LSTM-CRF model that is not re-trained on the combined training and validation data. FLDCRF-s outperforms this LSTM-CRF model (F1 score 85%) as well on the test data.

- *Model selection*:

FLDCRF-s models only need to be optimized for number of hidden layers N_h and number of hidden states N_s . Each LSTM model (LSTM and LSTM-CRF, with a given number of hidden units N_{hls}), on the other hand, need to be tuned for the optimum number of training epochs; thus requiring careful monitoring on validation performance in order to avoid overfitting. We analyze the selection process of the best performing models in each family, viz., FLDCRF-s and LSTM-CRF below.

Table 6 presents FLDCRF-s performance on validation data. The models continue to perform better on increasing N_s upto 6 with $N_h = 1$, but decline on increasing N_h , giving clear indications to stop testing more hyperparameters. On the other hand, LSTM-CRF model performances vary rapidly across the training epochs (see Table 7) for most of the N_{hls} settings. While in some cases (e.g., $N_{hls} = 10, 14$ etc.) model performance goes down rapidly after attaining the maximum, there are cases (e.g., $N_{hls} = 100, 200$

Table 5: F1 scores of the models for the experiment proposed by Wagner et al. (2014) on UCI gesture phase data. The best LSTM-CRF model on the validation set achieved a F1 score 84.88 without retraining on the training+validation set.

Model	F1-score (%)
CRF	81.35
LDCRF	88.60
FLDCRF-s	88.60
LSTM	83.42
LSTM-CRF	84.28 (84.88)
Multi-layered Perceptron (Wagner et al., 2014)	82.34

etc.) where the validation F1-score varies randomly across training epochs, without any discernible pattern against the hyperparameter. Such random variation demands very careful monitoring over the validation outcomes. Even when optimized for the training epochs for each N_{hls} setting (see column 7 of Table 7), F1 scores do not reveal any pattern over the N_{hls} 's. Such a behaviour additionally brings in the problem of considering widespread values of N_{hls} .

Table 6: FLDCRF-s validation performance on experiment 2 (Wagner et al., 2014) on UCI gesture phase data.

FLDCRF-s	1/1	1/2	1/3	1/4	1/5	1/6	1/7	Best	Worst	Std
F1	80.07	80.79	80.87	81.31	81.6	82.96	81.92	82.96	76.47	2.11
FLDCRF-s	2/1	2/2	2/3	2/4	3/1	3/2	3/3			
F1	77.57	77.84	78.97	78.03	76.84	76.60	76.47			

FLDCRF-s models do not vary rapidly (standard deviation 2.11) across the choices of the hyperparameters, and rather indicate a worst case F1-score of 76.47%. On the other hand, a high percentage (18%) of the LSTM-CRF models generate poor F1 performance (<75%) with a notably low worst case performance (62%). In addition, LSTM-CRF models display a considerably higher standard deviation (5.11) among validated models.

Since there is no rule to select the best performing hyperparameter(s) for a LSTM-CRF, large and rapid variation in model performance across hyperparameter settings, high-percentage of poor performances with notably low worst cases, and longer training times (see below) make it very difficult for a fast and reliable LSTM model selection. We will see similar LSTM behaviour in all our experiments.

- *Consistency:*

Table 7: LSTM-CRF validation performance on experiment 2 (Wagner et al., 2014) on UCI gesture phase data. N_{els} denotes number of training epochs and N_{hls} represents the number of hidden units in the model.

$N_{hls} \backslash N_{els}$	2	5	7	10	14	20	50	100	200	Best	Worst	Std
100	68.9	78.7	78.8	81.3	80	83.9	73.4	78.8	72.1	87.17	62.45 (18%<75)	5.59
200	75	81.1	82.5	71.9	75.6	76.5	81.4	68.2	86.4			
300	79.9	80.1	83.6	67.4	87.2	80.3	82.8	75.9	77.7			
400	81.8	78.6	83.4	63.3	80.2	79.1	77.3	75.3	84.7			
500	83.4	79.4	84.2	62.5	76.5	81.2	77.8	75.1	76.9			
Best	83.4	81.1	84.2	81.3	87.2	83.9	82.8	78.8	86.4			
Std	5.9	1.01	2.17	7.69	4.54	2.71	3.7	3.9	5.88			

Although FLDCRF-s offers easier model selection than LSTM-CRF, LSTM-CRF outperforms FLDCRF-s on the best validation performance, yielding a F1-score 87% compared to 83% by FLDCRF-s. However, this optimized LSTM-CRF model, when re-trained with same hyperparameters on training+validation data, fails to outperform FLDCRF-s on the test data (see Table 5), achieving 84.3% compared to 88.6% by FLDCRF-s. Furthermore, considering variable performance for LSTM-CRF on re-training the model, we apply the best (pre-trained) LSTM-CRF model (with 87% on validation) on the test data. However, this model (F1-score 84.9%) also fails to outperform FLDCRF-s on the test data (see Table 5).

In addition to a tedious model selection process, LSTM-CRF displays inconsistent performance across validation and test data. This inconsistency is reflected throughout the experiments in the paper.

- *Computation times:*

We compare training and inference times required by the FLDCRF and LSTM models in Fig. 11. While few FLDCRF-s models and LSTM required similar computation times for inference, computation time required to train the FLDCRF-s models is notably lower than the LSTM models. Furthermore, FLDCRF computation times can be significantly reduced by GPU implementation and careful optimization, allowing more hidden layers (N_h) and states (N_s) for large datasets, if necessary. We consider the GPU implementation and application of FLDCRF on big data as our future work.

6.1.2 EXPERIMENT 2: NESTED CV

Since the training and test datasets in Experiment 1 were small, we consider the entire UCI gesture phase dataset in this experiment. We perform a nested cross-validation. Since there are 7 large sequences in the dataset (see Section 4.1), we divide the data in 7 outer folds, one for testing each sequence. In each case, we select the models by a 6-fold (leave-one-out)

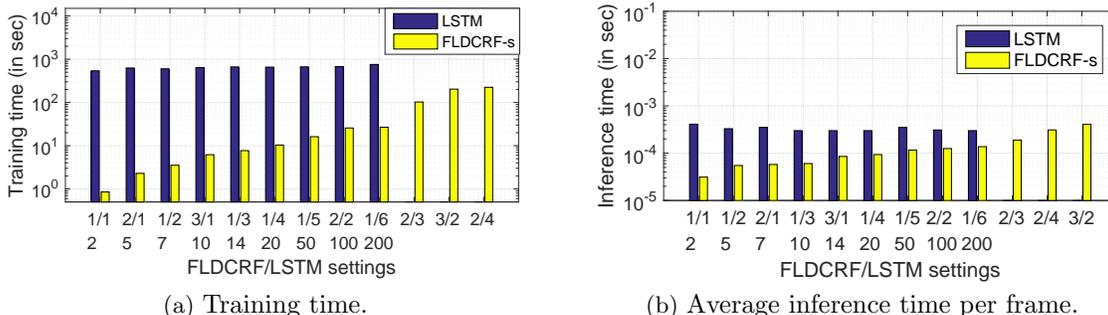


Figure 11: Training and inference times of different FLDCRF and LSTM models on UCI gesture phase data. Training time shown for each FLDCRF-s is until model convergence on training data, with the default criteria in Stan. Training time shown for LSTM is for 500 training epochs.

cross-validation on the remaining sequences.

- *Test performance:*

Table 8 compares different models on the test sets. FLDCRF-s boosts LDCRF performance on 4 sets (sets 1, 2, 3, 7), with notable improvements ($\sim 3\%$ and $\sim 2\%$) on sets 1 and 3. As a result, FLDCRF-s improves the overall LDCRF performance to outshine the overall LSTM (and LSTM-CRF) performance on the test sets. FLDCRF-s individually outperforms LSTM and LSTM-CRF models on 4 (out of 7) sets each.

Table 8: Different model performances on the nested CV experiment on UCI gesture phase data.

Model	F1 score							
	1	2	3	4	5	6	7	Average
LDCRF	85.6	76.92	87.81	92.6	92.76	91.83	76.1	86.23
FLDCRF-s	88.4	76.95	89.69	92.6	92.76	91.38	76.17	86.85
LSTM	90.95	81.61	82.59	94.41	92.51	89.12	73.5	86.39
LSTM-CRF	91.99	81.51	81.60	93.40	92.76	86.87	74.06	86.03

- *Model selection:*

Tables 9 and 10 present FLDCRF-s and LSTM cross-validation F1 outcomes on sets 4 and 7 respectively. Similar to experiment 1, LSTM models display rapid variation in performance across training epochs, demanding careful tracking of validation performance. Few such instances are highlighted in Table 10. FLDCRF-s does not require to tune the training epochs, and variation among adjacent hyperparameters (see Table 9) is neither rapid nor random, with observable decay in performance due to overfitting beyond certain settings of the hyperparameters.

Table 11 highlights the overall FLDCRF-s and LSTM cross-validation performance on the inner loops (for selecting models). We report the best, worst case performances and

Table 9: FLDCRF-s validation performance on nested validation sets 4 and 7 of the nested CV experiment on UCI gesture phase data.

Set 4									
FLDCRF	1/1	1/2	1/3	1/4	1/5	1/6	Best	Worst	Std
F1	86.38	85.58	86.73	86.5	86.33	85.96	86.73	85.01	0.50
FLDCRF	2/1	2/2	2/3	2/4	3/1	3/2			
F1	86.44	86.44	86.5	85.01	86.6	85.89			
Set 7									
FLDCRF	1/1	1/2	1/3	1/4	1/5	1/6	Best	Worst	Std
F1	88.3	89.05	89.19	88.7	88.52	88.96	89.23	88.3	0.28
FLDCRF	2/1	2/2	2/3	2/4	3/1	3/2			
F1	89.05	88.94	89.23	89.11	88.94	89.09			

Table 10: LSTM validation performance on nested validation sets 4 and 7 of the nested CV experiment on UCI gesture phase data. N_{els} denotes number of epochs trained and N_{hls} represents the number of hidden units in the model.

Set 4												
$N_{hls} \backslash N_{els}$	2	5	7	10	14	20	50	100	200	Best	Worst	Std
100	86.3	86.7	86.3	87.2	85.6	83.4	84.6	85.3	84.6	87.96	81.64	1.43
200	85.3	87.3	86.3	84.6	85.6	86.6	82.5	86.0	83.9			
300	85.0	87.8	85.5	87.9	85.3	87.5	85.5	86.6	85.5			
400	84.8	85.7	86.2	86.7	85.9	87.9	83.5	86.3	85.7			
500	85.5	86.0	84.7	86.1	85.3	83.1	83.5	81.6	86.4			
Std	0.59	0.71	0.69	1.23	0.28	2.31	1.14	2.05	0.99			
Set 7												
$N_{hls} \backslash N_{els}$	2	5	7	10	14	20	50	100	200	Best	Worst	Std
100	86.6	86.9	86.2	86.4	86.6	86.8	87.7	90.1	88.6	90.72	80.31	1.65
200	86.5	88.5	86.6	85.7	86.7	89.3	86.7	88.2	87.7			
300	88.3	89.0	87.2	86.8	85.9	90.1	86.2	88.3	87.8			
400	87.3	87.9	87.5	86.5	86.9	90.7	86.9	88.1	87.4			
500	85.7	88.7	88.2	87.2	87.6	80.3	87.0	87.1	88.2			
Std	0.99	0.83	0.78	0.56	0.59	4.27	0.55	1.08	0.47			

standard deviation among models on validation sets. We also report percentage of poor performing models, i.e., validated models (hyperparameters) below a given threshold (F1-

score 85%). FLDCRF-s produces the lower worst case performance with lesser standard deviation among models than LSTM in most of the sets.

High percentage of LSTM models performing below threshold across most of the sets is also noticeable and brings about difficulty to select the optimum models, by requiring careful monitoring and having to consider widespread hyperparameter settings. 5% of all considered FLDCRF-s models across 7 cross-validation sets report below F1-score 85, as compared to 17.1% for LSTM models. In set 7, although 2.2% of LSTM models report F1 below 85, a staggering 69% of the models perform below F1-score 88, while FLDCRF-s has a worst case F1 of 88.3.

Table 11: Summary of cross-validation performance by FLDCRF-s and LSTM on the inner loops of nested CV experiment on the UCI gesture phase data.

Model	F1-score (best)						
	1	2	3	4	5	6	7
FLDCRF-s	87.44	86.55	88.88	86.73	88.59	86.32	89.23
LSTM	88.24	87.32	89.85	87.96	86.71	87.27	90.72
Model	F1-score (worst)						
	1	2	3	4	5	6	7
FLDCRF-s	86.02	84.38	84.89	85.01	86.91	85.03	88.3
LSTM	84.24	82.17	85.17	81.64	82.24	78.58	80.31
Model	F1-score (std)						
	1	2	3	4	5	6	7
FLDCRF-s	0.38	0.83	1.29	0.55	0.47	0.50	0.27
LSTM	0.90	0.70	1.07	1.43	0.94	1.36	1.65
Model	# F1-score < 85 (% of all validated models)						
	1	2	3	4	5	6	7
FLDCRF-s	0	28.5	7.14	0	0	0	0
LSTM	8.9	20	0	26.7	37.8	24.4	2.2 ¹

- *Consistency:*

LSTM outperforms FLDCRF-s on the best cross-validation performances (see F1-score (best) in Table 11) on most (6 out of 7) of the sets. However, such optimized models fail to consistently outperform FLDCRF-s on the test sets (see Table 8), and manage to perform better only on 3 (out of 7) sets. The average performance reported by LSTM on the test data is also about 0.5% less than FLDCRF-s.

- *Computation times:*

Figure 12 compares considered FLDCRF-s and LSTM training times. LSTM training times are shown for 500 training epochs, while the FLDCRF-s training times are until convergence (by the default criteria for BFGS in Stan). Displayed results are quite similar

1. 68.9% of the LSTM models below F1-score 88, while FLDCRF-s reports worst case F1 of 88.3.

to experiment 1, only magnified by the amount of training data. As mentioned earlier, FLDCRF training times can be further reduced by GPU implementation.

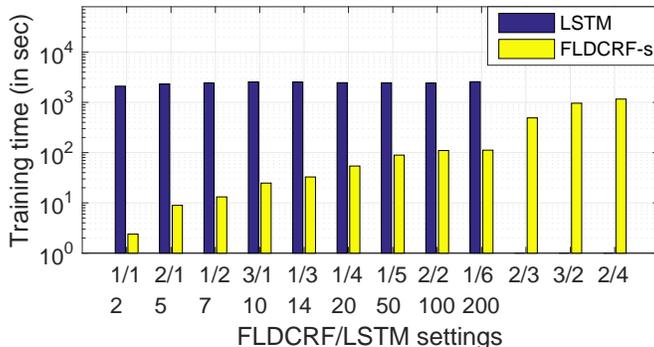


Figure 12: Average training time required by different FLDCRF settings per outer fold of the nested CV experiment on the UCI gesture phase data.

6.2 UCI Opportunity Dataset

We perform our multi-label sequence tagging experiment on the UCI Opportunity Dataset. We compare several single-label (CRF, LDCRF, FLDCRF-s, LSTM, LSTM-CRF) and multi-label (FCRF, CCRF, FLDCRF-m1, FLDCRF-m2, LSTM-m) models in the nested CV experiment.

6.2.1 EXPERIMENT 1: NESTED CV

This experiment has 5 outer loops, one for testing on each ADL sequence (see Section 4.2). In each case, we select the models by 4-fold cross-validation on the remaining 4 ADL sequences. Sequence details for this experiment are illustrated in Table 2. Tabulated data (features and labels) for this experiment is available here².

- *Test performance:*

We present individual labeling results (F1-scores) from different models in Tables 12 (for locomotion activity) and 13 (for HL activity). The results for the joint labeling task are presented in Table 14.

Similar to the nested CV experiment on the UCI gesture phase data, FLDCRF-s improves the overall LDCRF performance on the dataset and outperforms LSTM and all other models. LSTM with the regular softmax layer for classification outperforms other LSTM models and LDCRF. LDCRF does not significantly improve CRF performance by modeling the latent dynamics, neither do the multi-label models FCRF, CCRF by considering joint learning. However, FLDCRF-m2 achieves notable improvement ($\sim 0.4\%$ on average) over these models by considering latent-dynamic interactions (cotemporal and first-order Markov, see Section 5.2) among the label categories (locomotion and HL). It should be noted that, multiple hidden layers (1-3) for the same activity label were

2. <https://github.com/satyajitneogiju/FLDCRF-for-sequence-labeling>

involved in the validation process of FLDCRF-s. So, multiple hidden layers for each label category can be employed in a FLDCRF-m2 model in order to further improve performance. FLDCRF-m1 only includes co-temporal interaction among the hidden layers (see Section 5.2), and achieves marginal improvement over other CRF models.

Table 12: Different model performances on the ‘locomotion’ activity label during nested CV experiment on UCI opportunity data.

Model	F1 score (Locomotion)					
	1	2	3	4	5	Average
CRF	48.26	64.4	84.01	82.32	86.87	73.17
LDCRF	48.3	64.4	84.04	82.32	86.87	73.19
FLDCRF-s	48.3	64.4	84.04	85.83	86.87	73.89
FCRF	48.16	64.42	83.93	82.53	86.84	73.18
CCRF	48.16	64.32	83.95	82.56	86.74	73.15
FLDCRF-m1	48.23	64.4	83.98	82.84	86.87	73.26
FLDCRF-m2	48.23	64.52	84.23	83.88	86.93	73.56
LSTM	46.41	64.28	77.64	91.73	87.94	73.60
LSTM-CRF	46.02	59.98	73.48	89.19	81.16	69.97
LSTM-m	46.60	61.47	78.66	90.63	88.38	73.15

Table 13 presents the labeling performance of different models on continuous tagging of the high-level (HL) activity. LDCRF, FLDCRF-s, FCRF and CCRF produce similar overall performance, marginally improving over the simple LCCRF (or CRF) model. LSTM-CRF outperforms LSTM and LSTM-m models and achieves a similar overall performance to that of the LCCRF model. FLDCRF-m2, by considering the latent-dynamic interactions in the joint labeling task, significantly (by at least 0.75% on average) outperforms all CRF and LSTM models, with significant (>3%) improvement over LDCRF on set 1. As in case of locomotion, FLDCRF-m1 achieves marginal improvement over other CRF models.

Table 14 summarizes the joint labeling performance of different models. FLDCRF-s performs best among the single-label sequence models, while FLDCRF-m2 outperforms all single and multi-label sequence models. LSTM models fail to perform consistently across the test sets, and produces relatively poorer overall performance. The multi-label LSTM-m achieves best results among the LSTM models. FLDCRF-m2 improves the LDCRF, FCRF, CCRF and CRF performance on all sets, demonstrating the significance of modeling latent-dynamic interactions among different label categories. FLDCRF-m2 outshines the LSTM model on 3 (out of 5) sets, and outperforms LSTM-CRF and LSTM-m models on 4 sets each, with significantly outperforming (by >1.5%) all the models on the entire dataset.

- *Model selection:*

We present the validation performance of the best models from each family, FLDCRF-m (written for FLDCRF-m2) and LSTM-m, on cross-validation sets 2 and 5 in Tables

Table 13: Different model performances on the ‘high-level’ activity (HL) label during nested CV experiment on UCI opportunity data.

Model	F1 score (HL)					
	1	2	3	4	5	Average
CRF	80.29	97.43	99.57	100	100	95.46
LDCRF	81.04	97.57	99.57	100	100	95.64
FLDCRF-s	81.04	97.55	99.56	100	100	95.63
FCRF	80.92	97.7	99.6	100	100	95.64
CCRF	80.91	97.71	99.57	100	100	95.64
FLDCRF-m1	81.38	97.73	99.72	99.88	100	95.74
FLDCRF-m2	84.29	98.1	99.56	100	100	96.39
LSTM	94.26	90.57	100	100	89.04	94.77
LSTM-CRF	81.21	98.04	100	100	98.02	95.45
LSTM-m	80.05	97.01	99.57	100	98.28	94.98

Table 14: Overall model performances on the joint sequence labeling task during nested CV experiment on UCI opportunity data.

Model	F1 score (Overall)					
	1	2	3	4	5	Average
CRF	57.67	79.69	91.58	91.16	93.43	82.71
LDCRF	58.21	79.83	91.59	91.16	93.43	82.84
FLDCRF-s	58.21	79.80	91.58	92.92	93.45	83.19
FCRF	58.05	79.97	91.56	91.27	93.42	82.85
CCRF	58.05	79.93	91.55	91.28	93.37	82.84
FLDCRF-m1	58.42	79.98	91.71	91.3	93.43	82.97
FLDCRF-m2	60.53	80.4	91.67	91.94	93.47	83.60
LSTM	67.78	73.52	88.82	95.86	84.09	82.01
LSTM-CRF	57.19	78.07	86.74	94.60	88.64	81.05
LSTM-m	56.67	77.83	88.90	95.31	92.50	82.24

15 and 16 respectively. We omit the detailed results for other sets to save space, and summarize all the best, worst and standard deviation performances on the inner loops of nested CV in Table 17.

LSTM-m gives rapidly fluctuating validation performance across epochs for most of the considered N_{hls} setups (see Table 16), especially during cross-validation on set 5. It also produces significantly lower worst case performance with high percentage of models giving notably poor results (19.5% of all models across 5 sets with $F1 < 72$, see Table 17 for individual sets). FLDCRF-m, on the other hand, shows much discipline against the

choices of the hyperparameters, and shows decline in performance beyond certain N_h/N_s settings (see Tables 15 and 17), making it much easier to work with.

Table 15: FLDCRF-m2 validation performance on nested cross-validation sets 2 and 5 on UCI opportunity data.

Set 2									
FLDCRF	1/1	1/2	1/3	1/4	1/5	1/6	Best	Worst	Std
F1	76.17	76.18	76.23	76.2	76.63	76.2	76.84	75.34	0.33
FLDCRF	2/1	2/2	2/3	2/4	2/5	2/6			
F1	76.32	76.84	76.33	76.51	76.32	76.32			
FLDCRF	3/1	3/2	3/3						
F1	76.31	76.65	75.34						

Set 5									
FLDCRF	1/1	1/2	1/3	1/4	1/5	1/6	Best	Worst	Std
F1	73.27	74.48	73.29	73.27	74.37	73.62	74.48	72.78	0.46
FLDCRF	2/1	2/2	2/3	2/4	2/5	2/6			
F1	73.32	73.36	74.11	73.3	73.3	73.53			
FLDCRF	3/1	3/2	3/3						
F1	73.33	73.31	72.78						

Table 16: LSTM-m validation performance on nested cross-validation sets 2 and 5 on UCI opportunity data.

Set 2											
$N_{hls} \backslash N_{els}$	5	10	25	50	75	150	300	500	Best	Worst	Std
100	81.1	79.4	79.6	79.7	79.3	81.1	81.3	70.1	83.23	67.98	3.49
200	74.0	80.2	83.2	76.4	75.6	78.2	75.0	73.6			
300	80.5	77.6	81.8	75.5	80.3	78.2	80.8	72.5			
400	75.3	75.1	77.5	74.3	77.6	76.1	76.5	68.1			
500	77.6	75.9	77.2	76.2	74.7	76.1	75.9	67.9			

Set 5											
$N_{hls} \backslash N_{els}$	5	10	25	50	75	150	300	500	Best	Worst	Std
100	71.6	73.3	71.7	77.4	68.6	81.1	72.2	77.6	85.7	67.49	4.17
200	75.0	70.9	71.3	80.6	67.4	73.6	72.2	78.9			
300	78.3	71.1	70.1	71.4	70.0	70.1	69.2	79.6			
400	71.6	70.4	78.8	72.0	71.3	74	72.2	78.9			
500	72.1	71.4	73.2	72.3	78.8	85.7	70.8	79.6			

Table 17: Summary of cross-validation performance by FLDCRF-s and LSTM on the inner loops of nested CV experiment on the UCI opportunity data.

Model	F1 score (best)				
	1	2	3	4	5
FLDCRF-m2	87.54	76.84	75.12	74.31	74.48
LSTM-m	87.70	83.23	79.45	80.15	85.70
Model	F1 score (worst)				
	1	2	3	4	5
FLDCRF-m2	86.96	75.34	73.25	72.67	72.78
LSTM-m	79.26	67.98	68.81	65.23	67.49
Model	F1 score (std)				
	1	2	3	4	5
FLDCRF-m2	0.17	0.33	0.66	0.45	0.46
LSTM-m	2.00	3.49	3.35	4.26	4.17
Model	% F1 < 72				
	1	2	3	4	5
FLDCRF-m2	0	0	0	0	0
LSTM-m	0	10	17.5	27.5	42.5

- *Consistency:*

As in our earlier experiments, LSTM models (LSTM-m) are considerably better than FLDCRF-m to give the best validation results across all sets (see Table 17). However, such optimized LSTM-m models fail to consistently beat FLDCRF-m on the test sets, rather producing inconsistent performance across test sets (managing superior test performance only on 1 out of 5 sets, see Table 14). This is true for all other LSTM variants tested in this paper. Such inconsistency among validation and test performance raises serious concerns about practical applications of LSTM models.

- *Computation times:*

Figure 13 compares training and inference (to include) times required by different FLDCRF-m and LSTM-m models. As mentioned earlier, FLDCRF-m is run on a CPU and can be made faster by GPU implementation.

6.3 Multi-view Experiment

In this section, we examine the multi-view LDCRFs (Song et al., 2012) (see Fig. 6) on the UCI opportunity data. We perform nested CV experiments separately on the locomotion and HL activity labels of the UCI opportunity data. As described in Section 4.2, the feature set x_t for opportunity data comprises of 113 body-worn sensor features ($x_{1,t}$) and 32 shoe sensor ($x_{2,t}$) features i.e., $x_t = \{x_{1,t}, x_{2,t}\}$. We consider three different LDCRFs (with observations $x_{1,t}$, $x_{2,t}$ and x_t respectively, see Fig. 1c) and a coupled-linked MVLDCRF for

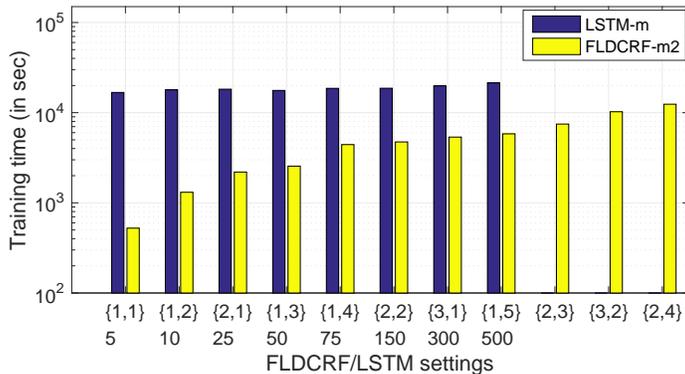


Figure 13: Average training time required by FLDCRF-m and LSTM-m per outer fold of UCI Opportunity data.

comparison. We also consider the FLDCRF-s performance with x_t . We present results of compared models on 10 test sets in Table 18.

LDCRF-body, LDCRF-shoe and LDCRF-(body+shoe) corresponds to LDCRFs with observations $x_{1,t}$, $x_{2,t}$ and x_t respectively. LDCRF-averaged gives the mean of LDCRF-body and LDCRF-shoe performances. MVLDCRF-averaged performance corresponds to the training and weighted inference mechanism by Song et al. (2012), while MVLDCRF-body layer and MVLDCRF-shoe layer gives the F1 scores obtained from individual layers of Fig. 6.

In most of the cases (sets 1, 3 of HL; set 5 of locomotion) where MVLDCRF-averaged (Song et al., 2012) outperformed the LDCRF with x_t , i.e., LDCRF-(body+shoe), either LDCRF-body (set 5 of locomotion) or LDCRF-shoe (sets 1, 3 of HL) has outperformed all other models, including MVLDCRF-averaged. In most other cases, one of the LDCRFs (LDCRF-(body+shoe), LDCRF-body, LDCRF-shoe) has outperformed MVLDCRF-averaged.

MVLDCRF-averaged marginally outperforms all 3 LDCRF models (i.e., with $x_{1,t}$, $x_{2,t}$ and x_t) only on set 2 of HL, where LDCRF-body and LDCRF-shoe performances are similar. We thus argue that an improvement by a MVLDCRF-averaged (Song et al., 2012) over all three LDCRFs (with $x_{1,t}$, $x_{2,t}$ and x_t) is possible only when the two individual LDCRFs with $x_{1,t}$ and $x_{2,t}$ perform similarly, as in set 2 of HL. It is advisable to test the three LDCRFs (or FLDCRF-s for improved performance) in case a feature distribution ($x_{1,t}$ and $x_{2,t}$) is available, else utilize the entire feature set x_t with a FLDCRF-s and let the model learn the different interacting latent dynamics within the features and labels. Figuring out the distributions ($x_{1,t}$, $x_{2,t}$ etc.) if any, where two LDCRFs perform similarly is very much the only way for a MVLDCRF to perform better than all three LDCRFs (and possibly FLDCRF-s). However such a task is quite tedious, if not impossible.

7. Discussion

We have shown difficulties to select LSTM models on validation data across our experiments in the paper. LSTM models, although producing some excellent performance on validation data for certain hyperparameter settings, fail to be consistent on the test data. FLDCRF-

Table 18: Comparing different models by distributing body-worn (113) and shoe sensor (32) features. LDCRF-body, LDCRF-shoe and LDCRF-(body+shoe) corresponds to $x_{1,t}$, $x_{2,t}$ and x_t respectively. LDCRF-averaged gives the mean of LDCRF-body and LDCRF-shoe outputs. MVLDCRF-averaged performance corresponds to the training and weighted inference mechanism by Song et al. (2012), while MVLDCRF-body layer and MVLDCRF-shoe layer gives the F1-scores obtained from individual layers of Fig. 6.

Model	HL					
	Set1	Set2	Set3	Set4	Set5	Average
LDCRF-body	79.96	97.13	99.57	99.83	100	95.29
LDCRF-shoe	98.48	97.83	99.81	99.97	99.94	99.21
LDCRF-averaged	89.22	97.48	99.69	99.9	99.97	97.25
LDCRF-(body+shoe)	81.04	97.57	99.57	100	100	95.64
FLDCRF-s(body+shoe)	81.04	97.55	99.56	100	100	95.63
MVLDCRF-body layer	86.63	97.96	99.63	100	100	96.84
MVLDCRF-shoe layer	91.81	98.18	99.71	100	100	97.94
MVLDCRF-averaged	89.22	98.07	99.67	100	100	97.39

Model	Locomotion					
	Set1	Set2	Set3	Set4	Set5	Average
LDCRF-body	46.6	60.33	81.5	82.9	87.22	71.71
LDCRF-shoe	51.55	71.97	76.81	81.34	77.81	71.9
LDCRF-averaged	49.07	66.15	79.155	82.12	82.51	71.80
LDCRF-(body+shoe)	48.3	64.4	84.04	82.32	86.87	73.19
FLDCRF-s(body+shoe)	48.3	64.4	84.04	85.83	86.87	73.89
MVLDCRF-body layer	48.26	63.93	83.98	81.95	87.09	73.04
MVLDCRF-shoe layer	45.92	63.58	83.98	82.05	86.99	72.50
MVLDCRF-averaged	47.09	63.75	83.98	82	87.04	72.77

CRF outperforms LSTM on test data across all our experiments. Additionally, FLDCRF requires less computation times for training and inference, even without GPU implementation. Moreover, FLDCRF being a graphical model, it is very easy to include known dependency information among the variables for improved modeling. Such inclusion in

LSTM is less straightforward. Motivation behind FLDCRF, its graph structure and the mathematical model all are lucidly defined and are derived from each other. Although, we agree on the motivation and the mathematical formulations of LSTM, the modeling process is not very lucid and hard to comprehend. We summarize several validated modeling attributes of FLDCRF and LSTM in Table 19, in terms of user convenience.

With difficulties in model selection, inconsistency across validation and test data, longer training and average performance on test data, LSTM has concerns about its stability and reliability on practical deployment. On the other hand, FLDCRF models bring peace of mind with lucid intuition, ease of model selection, consistent and superior performance on test data with shorter training.

FLDCRF subsumes major sequential CRF variants, viz., LCCRF, LDCRF and DCRF (see Section 3) and outperforms each of these state-of-the-art models across experiments in the paper. FLDCRF performance can be further improved by considering more generic variants (utilize past features x_{t-1} , x_{t-2} etc. during modeling, second order Markov dependency along hidden layers, multiple hidden layers for each label category etc.). We look forward to GPU implementation of FLDCRF and apply on several sequence labeling tasks of computer vision and NLP.

Table 19: Comparing FLDCRF and LSTM modeling attributes. HP stands for hyperparameter. '++' stands for very convenient for user, '+' stands for convenient, '-' stands for inconvenient and '- -' stands for very inconvenient.

Attribute	Model	
	FLDCRF	LSTM
Test performance	+	-
Ease of model selection	++	- -
• Less types of HP	• +	• -
• Rule to select HP	• -	• -
• Tune training epochs	• +	• -
• HP starting point	• +	• -
• HP ending point	• +	• -
• Pattern among HP	• ++	• - -
• Worst validation case	• +	• - -
• % of poor models	• +	• -
• Variance among models	• +	• -
• Best validation case	• -	• +
Consistency across validation and test	++	- -
Computation time	++	- -
Including known dependency	++	- -
Intuitive	++	- -
Practical deployment	++	- -

8. Conclusion

We proposed FLDCRF, a single and multi-label generalization of LDCRF. We presented 3 different FLDCRF variants for single-label (FLDCRF-s), multi-label (FLDCRF-m) and multi-agent (FLDCRF-i) sequence prediction/tagging tasks. FLDCRF-s allows multiple interacting latent dynamics of the class labels and extends the capability of LDCRF, as well as outperforms LSTM and LSTM-CRF across multiple datasets. FLDCRF-m introduces hidden variables in a DCRF to accommodate latent dynamic interactions among different label categories, thereby improving DCRF performance and outperforming all state-of-the-art sequence models including CRF, LDCRF, FCRF, LSTM, LSTM-m on a joint sequence labeling task. We also described the LSTM model selection difficulties and its inconsistent performance across validation and test data (summarized in Table 19). By contrast, FLDCRF presents easier model selection, provides consistency across validation and test data and guarantees good performance on random model selection. FLDCRF also offers lucid model intuition and user flexibility over approach. We look forward to GPU implementation of FLDCRF and compare FLDCRF and LSTM on larger sequence datasets. Another interesting topic for future research is to apply end-to-end models with FLDCRF succeeded by CNN layers on popular computer vision problems like action recognition. We are also exploring the multi-agent FLDCRF-i model for joint intention prediction of pedestrians on crossing/not-crossing before Autonomous Vehicles. It is possible to extend the idea of factorized latent space in FLDCRF to heterogeneous (discrete and continuous) state space models.

Acknowledgments

We would like to thank Dr. Michael Hoy for his contributions in this work. This research is partially supported by the ST Engineering-NTU Corporate Lab through the NRF corporate lab@university scheme.

Appendix A. Latent Dynamic Conditional Random Fields (LDCRF)

In this section, we will describe the LDCRF (Morency et al., 2007) mathematical model.

A.1 Model

The task is to learn a probabilistic mapping between a time-series sequence of observed input features $\mathbf{x} = \{x_1, x_2, \dots, x_T\}$ and a sequence of observed classification labels $\mathbf{y} = \{y_1, y_2, \dots, y_T\}$. $y_t \in \Upsilon$, $\forall j = 1, 2, \dots, T$, where Υ is the set of classification labels. $\mathbf{h} = \{h_1, h_2, \dots, h_T\}$ denotes the hidden layer for capturing intrinsic dynamics within each label. Each label $\ell \in \Upsilon$ is associated with a set of hidden states \mathcal{H}_ℓ . \mathcal{H} is the set of all possible hidden states written as $\mathcal{H} = \bigcup_\ell \mathcal{H}_\ell$. \mathcal{H}_ℓ are disjoint $\forall \ell \in \Upsilon$. Each h_t is restricted to belong to the set \mathcal{H}_{y_t} , i.e., $h_t \in \mathcal{H}_{y_t}, \forall t = 1, 2, \dots, T$. The conditional model is defined as:

$$P(\mathbf{y} | \mathbf{x}, \theta) = \sum_{\mathbf{h}} P(\mathbf{y} | \mathbf{x}, \mathbf{h}, \theta) P(\mathbf{h} | \mathbf{x}, \theta). \quad (18)$$

Equation (18) can be re-written using the graph structure in Fig. 1c as:

$$\begin{aligned} P(\mathbf{y} | \mathbf{x}, \theta) &= \sum_{\mathbf{h}: \forall h_t \in \mathcal{H}_{y_t}} P(\mathbf{y} | \mathbf{h}, \theta) P(\mathbf{h} | \mathbf{x}, \theta) \\ &+ \sum_{\mathbf{h}: \exists h_t \notin \mathcal{H}_{y_t}} P(\mathbf{y} | \mathbf{h}, \theta) P(\mathbf{h} | \mathbf{x}, \theta). \end{aligned} \quad (19)$$

Applying model constraints, we can write the following:

$$P(y_t = \ell | h_t) = \begin{cases} 1, & h_t \in \mathcal{H}_{y_t=\ell} \\ 0, & h_t \notin \mathcal{H}_{y_t=\ell}. \end{cases} \quad (20)$$

The model in equation (19) can be simplified using equation (20) as:

$$P(\mathbf{y} | \mathbf{x}, \theta) = \sum_{\mathbf{h}: \forall h_t \in \mathcal{H}_{y_t}} P(\mathbf{h} | \mathbf{x}, \theta). \quad (21)$$

$P(\mathbf{h} | \mathbf{x}, \theta)$ is described using Conditional Random Field formulation given by:

$$P(\mathbf{h} | \mathbf{x}, \theta) = \frac{1}{\mathbf{Z}(\mathbf{x}, \theta)} \exp \left(\sum_k \theta_k \cdot F_k(\mathbf{h}, \mathbf{x}) \right), \quad (22)$$

where index k ranges over all parameters $\theta = \{\theta_k\}$ and $\mathbf{Z}(\mathbf{x}, \theta)$ is the partition function defined as:

$$\mathbf{Z}(\mathbf{x}, \theta) = \sum_{\mathbf{h}} \exp \left(\sum_k \theta_k \cdot F_k(\mathbf{h}, \mathbf{x}) \right). \quad (23)$$

The feature functions F_k 's are defined as:

$$F_k(\mathbf{h}, \mathbf{x}) = \sum_{t=1}^T f_k(h_{t-1}, h_t, \mathbf{x}, t),$$

Feature functions $f_k(h_{t-1}, h_t, \mathbf{x}, t)$ can be either an *observation* (also called *state*) function $s_k(h_t, \mathbf{x}, t)$ or a *transition* function $t_k(h_{t-1}, h_t, \mathbf{x}, t)$.

A.2 Training Model Parameters

Parameters of the model can be estimated by maximizing the conditional log-likelihood of the training data given by equation (24):

$$\mathbf{L}(\theta) = \sum_{n=1}^N \log P(\mathbf{y}^{(n)} | \mathbf{x}^{(n)}, \theta) - \frac{\|\theta\|^2}{2\sigma^2}, \quad (24)$$

where N is the total number of available labeled sequences. The second term in equation (24) is the log of a Gaussian prior with variance σ^2 .

A.3 Inference

Given a new test sequence \mathbf{x} , the inference task is given by:

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} P(\mathbf{y} | \mathbf{x}, \hat{\theta}), \quad (25)$$

Using model constraints, equation (25) can be re-written as:

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \sum_{\mathbf{h}: \forall h_t \in \mathcal{H}_{y_t}} P(\mathbf{h} | \mathbf{x}, \hat{\theta}), \quad (26)$$

We apply forward recursions of belief propagation for our inference as the problem of intention prediction must be solved online. In other words, at each time instant t , we compute the marginals $P(h_t | x_{1:t}, \theta)$ and sum them according to the disjoint sets of hidden states to obtain $P(y_t | x_{1:t}, \theta) = \sum_{h_t \in \mathcal{H}_{y_t}} P(h_t | x_{1:t}, \theta)$, $t = 1, 2, \dots$. Then, we infer the label y_t corresponding to the maximum probability. Forward-backward algorithm Rabiner (1989) and Viterbi algorithm Forney (1973) can also be applied for problems where online inference is not required.

References

- F. Sha and F. Pereira. "Shallow parsing with conditional random fields". In *Conference on Human Language Technology and North American Association for Computational Linguistics (HLT-NAACL)*, pages 213-220, 2003.
- L. R. Rabiner. "A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition". In *Proceedings of the IEEE*, Volume: 77, Issue: 2, Feb 1989, DOI: 10.1109/5.18626.

- J. Lafferty, A. McCallum, F. C. N. Pereira, “Conditional Random Fields: Probabilistic Models for Segmenting and Labeling Sequence Data”. In *Proceedings of the Eighteenth International Conference on Machine Learning*, Pages 282-289, June 2001.
- C. Sutton, K. Rohanimanesh, A. McCallum, “Dynamic Conditional Random Fields: Factorized Probabilistic Models for Labeling and Segmenting Sequence Data”. *Journal of Machine Learning Research*, 8:693-723, 2007.
- L. P. Morency, A. Quattoni, T. Darrell, “Latent-Dynamic Discriminative Models for Continuous Gesture Recognition”. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, IEEE Computer Society, 2007.
- G. Lample, M. Ballesteros, S. Subramanian, K. Kawakami, and C. Dyer, “Neural architectures for named entity recognition”. *arXiv:1603.01360*. [Online]. Available: <https://arxiv.org/abs/1603.01360>, March 2016.
- F. J. Ordez, D. Roggen, “Deep Convolutional and LSTM Recurrent Neural Networks for Multimodal Wearable Activity Recognition”. *Sensors* 16(1), 115; <https://doi.org/10.3390/s16010115>, 2016.
- R. Collobert, J. Weston, L. Bottou, M. Karlen, K. Kavukcuoglu, and P. Kuksa. “Natural language processing (almost) from scratch”. *Journal of Machine Learning Research*, 12:24932537, 2011.
- K. P. Murphy, “Dynamic bayesian networks: representation, inference and learning”. *Doctoral Dissertation*, University of California, Berkeley, 2002.
- A. McCallum, D. Freitag, F. C. N. Pereira, “Maximum Entropy Markov Models for Information Extraction and Segmentation”. In *Proceedings of the Seventeenth International Conference on Machine Learning*, Pages 591-598, 2000.
- Z. Huang, W. Xu, and K. Yu, “Bidirectional LSTM-CRF models for sequence tagging”. CoRR, abs/1508.01991, 2015.
- S. B. Wang, A. Quattoni, L. P. Morency, D. Demirdjian, T. Darrell, “Hidden Conditional Random Fields for Gesture Recognition”. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, Volume: 29, Issue: 10, Oct. 2007.
- X. Sun, L. P. Morency, D. Okanohara, J. Tsujii, “Modeling latent-dynamic in shallow parsing: a latent conditional model with improved inference”. In *Proceedings of the 22nd International Conference on Computational Linguistics*, Volume 1, Pages 841-848, 2008.
- Z. Ghahramani, M. I. Jordan, “Factorial Hidden Markov Models”. *Machine Learning*, 29: 245. <https://doi.org/10.1023/A:1007425814087>, 1997.
- S. Neogi, M. Hoy, W. Chaoqun, J. Dauwels, “Context Based Pedestrian Intention Prediction Using Factored Latent Dynamic Conditional Random Fields”. In *Proceedings of the 2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, DOI: 10.1109/SSCI.2017.8280970, 2017.

- D. Thai, S. H. Ramesh, S. Murty, L. Vilnis, A. McCallum, “Embedded-State Latent Conditional Random Fields for Sequence Labeling”, In *Proceedings of CoNLL 2018*, 2018.
- S. Hochreiter and J. Schmidhuber. “Long short-term memory”. *Neural computation*, 9(8):1735-1780, 1997.
- A. Ullah, J. Ahmad, K. Muhammad, M. Sajjad, S. W. Baik, “Action Recognition in Video Sequences using Deep Bi-directional LSTM with CNN features”. *IEEE Access*, Volume 6, Pages 841-848, DOI: 10.1109/ACCESS.2017.2778011, 2017.
- H. Zhu, I. Ch. Paschalidis, A. M. Tahmasebi, ”Context-based bidirectional-LSTM model for sequence labeling in clinical reports”. In *Proceeding of SPIE 10954, Medical Imaging 2019: Imaging Informatics for Healthcare, Research, and Applications*, DOI: 10.1117/12.2512103, March 2019.
- R. Caruana. “Multitask Learning”. *Machine Learning*, 28(1):4175, 1997.
- J. Wang, Y. Yang, J. Mao, Z. Huang, C. Huang, W. Xu, “CNN-RNN: A Unified Framework for Multi-label Image Classification”. *arXiv:1604.04573*, 2016.
- Y. Wei, W. Xia, M. Lin, J. Huang, B. Ni ; J. Dong, Y. Zhao, S. Yan, “HCP: A Flexible CNN Framework for Multi-Label Image Classification”. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Volume: 38 , Issue: 9, Pages 1901 - 1907, 2015.
- M. J. Choi, A. Torralba, A. S. Willsky, “A Tree-Based Context Model for Object Recognition”. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, Volume: 34 , Issue: 2, Pages 240-252, 2012.
- B. X. Nie, C. Xiong and S. C. Zhu, “Joint Action Recognition and Pose Estimation From Video”. In *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, DOI: 10.1109/CVPR.2015.7298734, 2015.
- S. Changpinyo, H. Hu, F. Sha, “Multi-Task Learning for Sequence Tagging: An Empirical Study”. *arXiv:1808.04151*, 2018.
- V. Pahuja, , A. Laha, , S Mirkin, , V. Raykar, L. Kotlerman, G. Lev, “Joint Learning of Correlated Sequence Labeling Tasks Using Bidirectional Recurrent Neural Networks”. In *Proceedings of the INTERSPEECH 2017*, 2017.
- Z. Li, J. Chao, M. Zhang, W. Chen, “Coupled Sequence Labeling on Heterogeneous Annotations: POS Tagging as a Case Study”. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, DOI: 10.3115/v1/P15-1172, 2015.
- G. Luo, X. Huang, C. Y. Lin, Z Nie, “Joint Named Entity Recognition and Disambiguation”. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 879888, 2015.
- M. Dredze, P. P. Talukdar, K. Crammer, “Sequence Learning from Data with Multiple Labels”. In *Proceedings of the ECML-PKDD 2009 workshop on learning from multi-label data*, MLD, 2009.

- Y. Shi, M. Wang, “A Dual-layer CRFs Based Joint Decoding Method for Cascaded Segmentation and Labeling Tasks”. In *Proceedings of the International Joint Conference on Artificial Intelligence 2007*, 2007.
- E. F. T. K. Sang, S. Buchholz, “Introduction to the CoNLL-2000 Shared Task: Chunking”. In *Proceedings of CoNLL-2000 and LLL-2000*, pages 127-132, Lisbon, Portugal, 2000.
- E. F. T. K. Sang, F. D. Meulder, “Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition”. In *CONLL '03 Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003*, Volume 4, Pages 142-147, Edmonton, Canada, 2003.
- A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, S. Savarese, “Social LSTM: Human Trajectory Prediction in Crowded Spaces”. In *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- R. Chavarriaga, H. Sagha, A. Calatroni, S. Digumarti, G. Trster, J. D. R. Milln, D Roggen. “The Opportunity challenge: A benchmark database for on-body sensor-based activity recognition”. *Pattern Recognition Letters*, 2013.
- Y. Song, L. P. Morency, R. Davis, “Multi-View Latent Variable Discriminative Models For Action Recognition”. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR) 2012*, 2012.
- W. Lu and H. T. Ng, “Better Punctuation Prediction with Dynamic Conditional Random Fields”. In *Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing*, pages 177186, 2010.
- R. C. B. Madeo, C. A. M. Lima, S. M. PERES, “Gesture Unit Segmentation using Support Vector Machines: Segmenting Gestures from Rest Positions”. In *Proceedings of the 28th Annual ACM Symposium on Applied Computing (SAC)*, p. 46-52, 2013.
- G. D. Forney Jr. “The Viterbi algorithm”. In *Proceedings of the IEEE*. 61 (3): 268278. doi:10.1109/PROC.1973.9030, March, 1973.
- P. K. Wagner, S. M. Peres, C. A. M. Lima, F. A. Freitas, R. C. B. Madeo, “Gesture Unit Segmentation Using Spatial-Temporal Information and Machine Learning”. In *Proceedings of the Twenty-Seventh International Florida Artificial Intelligence Research Society Conference*, 2014.
- Thomas W, “Thomas’ answer to the question ‘How to select number of hidden layers and number of memory cells in an LSTM?’ posted at Artificial Intelligence Stack Exchange website”. URL: <https://ai.stackexchange.com/questions/3156/how-to-select-number-of-hidden-layers-and-number-of-memory-cells-in-an-lstm>, 2017.
- Multiple users, “Answers to the question ‘How to choose the number of hidden layers and nodes in a feedforward neural network?’ posted at Stack Exchange website”. URL: <https://stats.stackexchange.com/questions/181/how-to-choose-the-number-of-hidden-layers-and-nodes-in-a-feedforward-neural-netw/1097#1097>, 2018-2019.

- K. Eckhardt, “Choosing the right Hyperparameters for a simple LSTM using Keras”, URL: <https://towardsdatascience.com/choosing-the-right-hyperparameters-for-a-simple-lstm-using-keras-f8e9ed76f046>, 2018.
- B. Carpenter, D. Lee, M. A. Brubaker, A. Riddell, A. Gelman, B. Goodrich, J. Guo, M. Hoffman, M. Betancourt, P. Li, “Stan: A Probabilistic Programming Language”. *Journal of Statistical Software*, Volume VV, Issue II, Nov. 2017.
- MATLAB 8.0 and Statistics Toolbox 8.1, The MathWorks, Inc., Natick, Massachusetts, United States.
- M. Wollmer, M. Kaiser, F. Eyben, B. Schuller, G. Rigoll, “LSTM-Modeling of continuous emotions in an audiovisual affect recognition framework”. *Image and Vision Computing*, 31 (2013) 153163, 2013.
- J. Huang, W. Zhou, Q. Zhang, H. Li, W. Li, “Video-Based Sign Language Recognition without Temporal Segmentation”. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, 2018.
- A. Zadeh, P. P. Liang, S. Poria, Prateek Vij, E. Cambria, L. P. Morency, “Multi-Attention Recurrent Network for Human Communication Comprehension”. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18)*, 2018.
- K. Greff, R. K. Srivastava, J. Koutnk, B. R. Steunebrink, J. Schmidhuber, “LSTM: A Search Space Odyssey”. In *IEEE TRANSACTIONS ON NEURAL NETWORKS AND LEARNING SYSTEMS*, Volume: 28, Issue: 10, Pages: 2222 - 2232, DOI: 10.1109/TNNLS.2016.2582924, Oct. 2017.
- S. Neogi, M. Hoy, K. Dang, H. Yu, J. Dauwels, “Context Model for Pedestrian Intention Prediction using Factored Latent-Dynamic Conditional Random Fields”. <https://arxiv.org/abs/1907.11881>, 2019.