

Scaling Matters in Deep Structured Prediction Models

Alexander Shevchenko, Anton Osokin

Structured Prediction

Dataset: $\mathcal{D} = \{(X^j, Y^j)\}_{j=1}^N$

Structure: $Y = \{y_1, \dots, y_L\}$

y_i highly correlated!

Tasks: sequence tagging, semantic segmentation, human pose estimation

Structured Prediction

Dataset: $\mathcal{D} = \{(X^j, Y^j)\}_{j=1}^N$

Structure: $Y = \{y_1, \dots, y_L\}$

y_i highly correlated!

Tasks: sequence tagging, semantic segmentation, human pose estimation

How to encode structure?

Energy-based Models

Score (negative energy): $F(Y \mid \Psi(X))$

Energy-based Models

Score (negative energy): $F(Y \mid \Psi(X))$

Associated likelihood: $p(Y \mid X) \propto \exp\{F(Y \mid \Psi(X))\}$

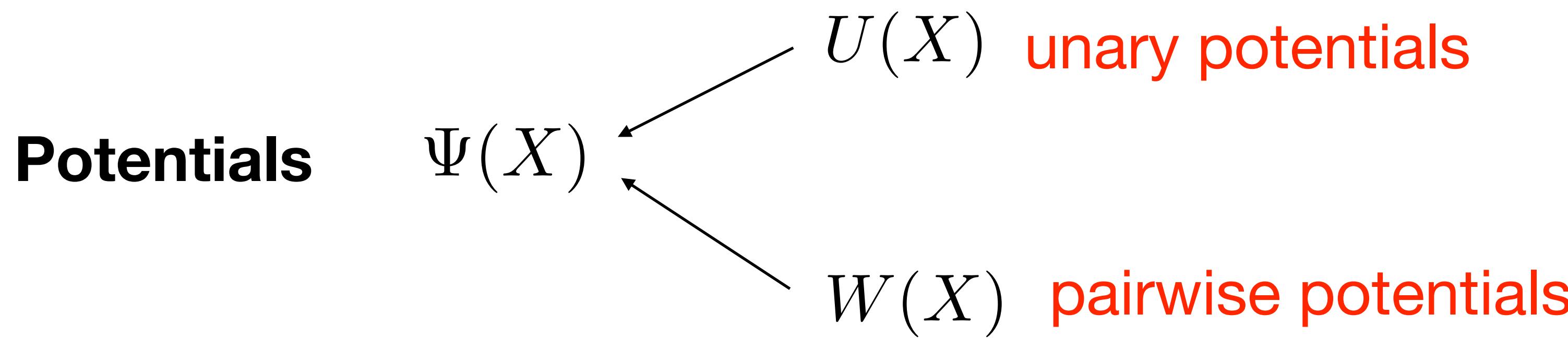
$$Z(X) = \sum_{Y'} \exp\{F(Y' \mid \Psi(X))\}$$

Energy-based Models

Score (negative energy): $F(Y \mid \Psi(X))$

Associated likelihood: $p(Y \mid X) \propto \exp\{F(Y \mid \Psi(X))\}$

$$Z(X) = \sum_{Y'} \exp\{F(Y' \mid \Psi(X))\}$$



Energy-based Models Training

Consider the parameterization of potentials Ψ_θ

Training objective $\frac{1}{N} \sum_{j=1}^N \mathcal{L}(X^j, Y^j, F_\theta(\cdot)) \rightarrow \min_\theta$

Energy-based Models Training

Consider the parameterization of potentials Ψ_θ

Training objective $\frac{1}{N} \sum_{j=1}^N \mathcal{L}(X^j, Y^j, F_\theta(\cdot)) \rightarrow \min_\theta$

Stage training

- + with proper tuning works almost always
- time-consuming
- not stable due to phase switching

Energy-based Models Training

Consider the parameterization of potentials Ψ_θ

Training objective
$$\frac{1}{N} \sum_{j=1}^N \mathcal{L}(X^j, Y^j, F_\theta(\cdot)) \rightarrow \min_{\theta}$$

Stage training

- + with proper tuning works almost always
- time-consuming
- not stable due to phase switching

Joint training

- increased optimization complexity
- + significantly faster
- + often results in better model

Energy-based Models Training

Consider the parameterization of potentials Ψ_θ

Training objective
$$\frac{1}{N} \sum_{j=1}^N \mathcal{L}(X^j, Y^j, F_\theta(\cdot)) \rightarrow \min_{\theta}$$

Stage training

- + with proper tuning works almost always
- time-consuming
- not stable due to phase switching

Joint training

- increased optimization complexity
- + significantly faster
- + often results in better model

Reason to fail: improper relative scaling

Motivation: OCR case study

OBJECTIVE	STAGE	JOINT
Cross-entropy	97.18 ± 0.12	96.48 ± 0.23
Structured SVM	96.97 ± 0.27	96.43 ± 0.39
Log-likelihood	97.15 ± 0.13	96.46 ± 0.25

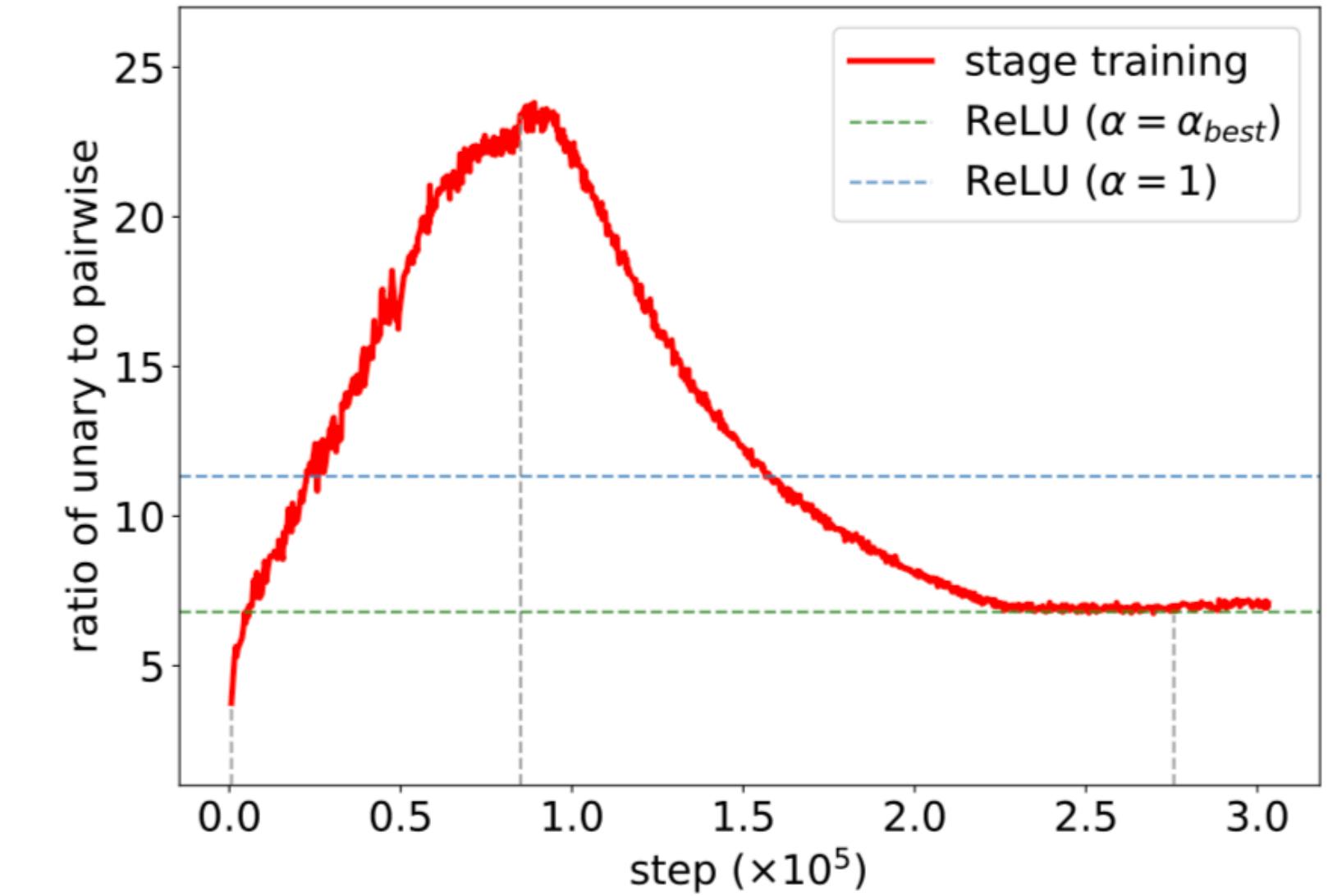
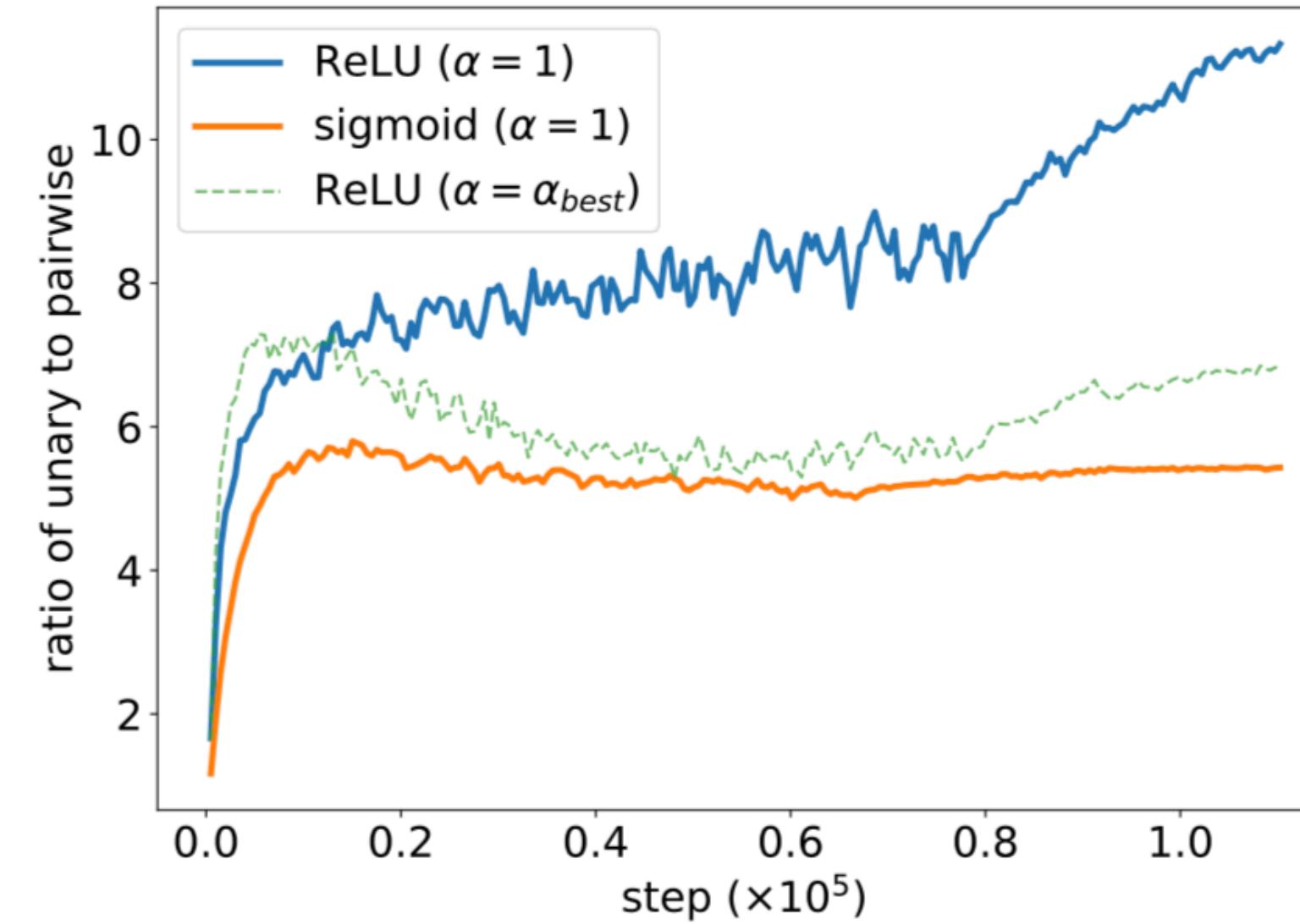
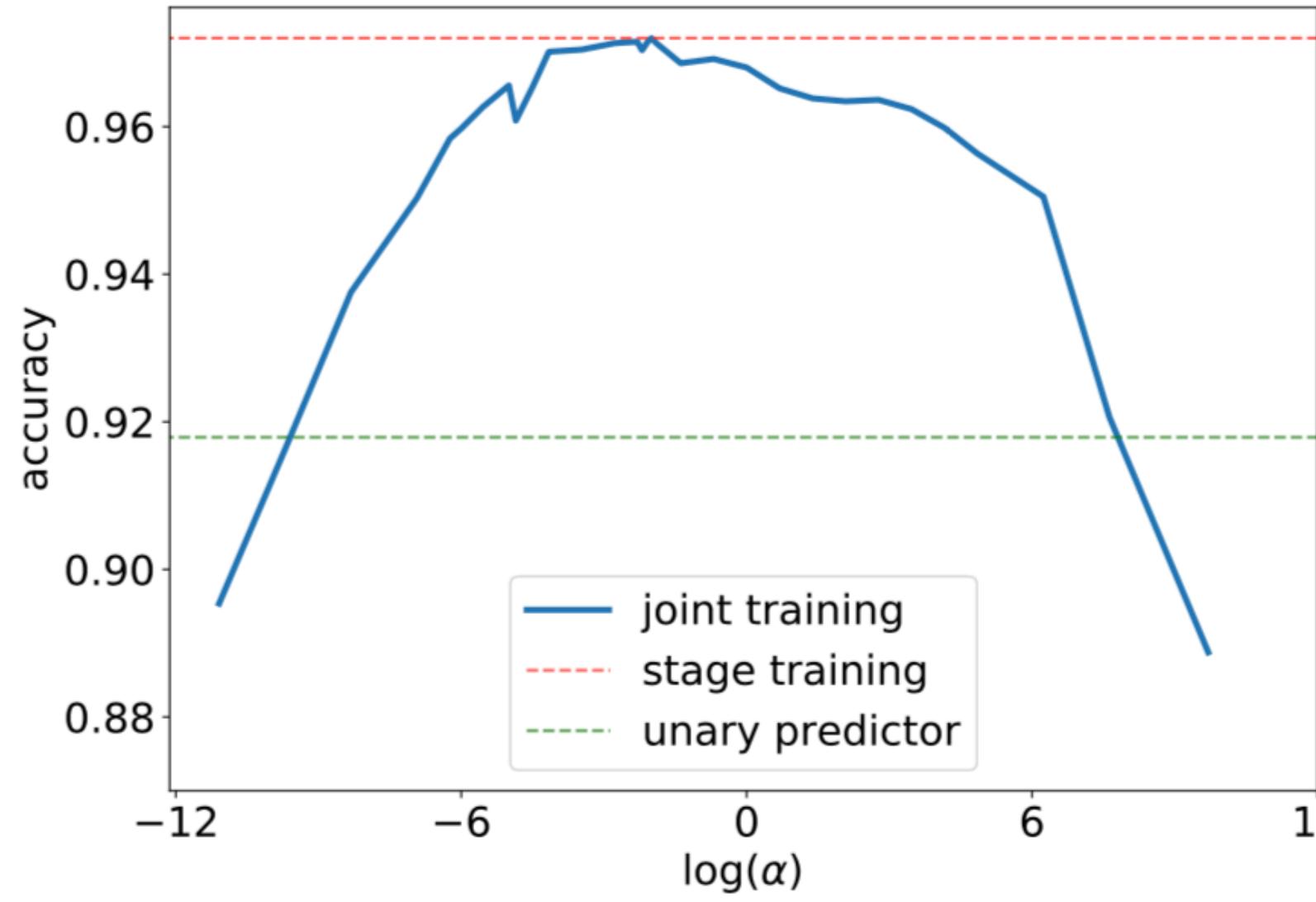
Motivation: OCR case study

Rescale unary

$$\hat{U}_\theta(X) = \alpha \cdot U(X)$$

OBJECTIVE	STAGE	JOINT
Cross-entropy	97.18 \pm 0.12	96.48 ± 0.23
Structured SVM	96.97 \pm 0.27	96.43 ± 0.39
Log-likelihood	97.15 \pm 0.13	96.46 ± 0.25

Motivation: OCR case study



Rescale unary

$$\hat{U}_\theta(X) = \alpha \cdot U(X)$$

OBJECTIVE	STAGE	JOINT
Cross-entropy	97.18 \pm 0.12	96.48 ± 0.23
Structured SVM	96.97 \pm 0.27	96.43 ± 0.39
Log-likelihood	97.15 \pm 0.13	96.46 ± 0.25

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y \mid X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y \mid X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Likelihood and marginals

$$P(Y \mid X, \theta) = \frac{\exp\{F(Y \mid X, \theta)\}}{Z(\theta)}$$

$$p_i(y_i \mid X, \theta) = \sum_{Y \setminus y_i} P(Y \mid X, \theta)$$

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y \mid X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Likelihood and marginals

$$P(Y \mid X, \theta) = \frac{\exp\{F(Y \mid X, \theta)\}}{Z(\theta)}$$

$$p_i(y_i \mid X, \theta) = \sum_{Y \setminus y_i} P(Y \mid X, \theta)$$

Inference

- ◆ Partition function, marginals (sum-product BP)
- ◆ MAP estimate (max-sum BP)

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y \mid X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Training objectives

$$\mathcal{L}_{\text{MLE}}(X, Y \mid \theta) = -\log P(Y \mid X, \theta)$$

Likelihood and marginals

$$P(Y \mid X, \theta) = \frac{\exp\{F(Y \mid X, \theta)\}}{Z(\theta)}$$

$$p_i(y_i \mid X, \theta) = \sum_{Y \setminus y_i} P(Y \mid X, \theta)$$

Inference

- ◆ Partition function, marginals (sum-product BP)
- ◆ MAP estimate (max-sum BP)

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y | X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Likelihood and marginals

$$P(Y | X, \theta) = \frac{\exp\{F(Y | X, \theta)\}}{Z(\theta)}$$

$$p_i(y_i | X, \theta) = \sum_{Y \setminus y_i} P(Y | X, \theta)$$

Inference

- ◆ Partition function, marginals (sum-product BP)
- ◆ MAP estimate (max-sum BP)

Training objectives

$$\mathcal{L}_{\text{MLE}}(X, Y | \theta) = -\log P(Y | X, \theta)$$

$$\mathcal{L}_{\text{cross}}(X, Y | \theta) = -\frac{1}{L} \sum_{i=1}^L \log p_i(y_i | X, \theta)$$

Linear-Chain CRF (LCCRF)

LCCRF score function

$$F(Y \mid X, \theta) = \sum_{i=1}^L U_\theta(x_i, y_i) + \sum_{i=1}^{L-1} W_\theta(y_i, y_{i+1})$$

Likelihood and marginals

$$P(Y \mid X, \theta) = \frac{\exp\{F(Y \mid X, \theta)\}}{Z(\theta)}$$

$$p_i(y_i \mid X, \theta) = \sum_{Y \setminus y_i} P(Y \mid X, \theta)$$

Inference

- ◆ Partition function, marginals (sum-product BP)
- ◆ MAP estimate (max-sum BP)

Training objectives

$$\mathcal{L}_{\text{MLE}}(X, Y \mid \theta) = -\log P(Y \mid X, \theta)$$

$$\mathcal{L}_{\text{cross}}(X, Y \mid \theta) = -\frac{1}{L} \sum_{i=1}^L \log p_i(y_i \mid X, \theta)$$

$$\begin{aligned} \mathcal{L}_{\text{SSVM}}(X, Y \mid \theta) = & \max_{Y'} \{F(Y' \mid X, \theta) + \Delta(Y', Y)\} \\ & - F(Y \mid X, \theta) \end{aligned}$$

Hamming distance

$$\Delta(Y', Y) = \frac{1}{L} \sum_{i=1}^L \mathbb{I}[y'_i \neq y_i]$$

Gaussian CRF (GCRF)

The energy function has different nature

GCRF score function

$$F(s \mid X, \theta) = -\frac{1}{2} s^T (W_\theta(X) + \lambda I) s + U_\theta(x) s$$

$$s = [s_{1,1}, \dots, s_{1,L} \mid \dots \mid s_{M,1}, \dots, s_{M,L}]$$

$$U_\theta(X) = [u_\theta^1(X) \mid \dots \mid u_\theta^M(X)] \quad W_\theta(X) \in \mathbb{R}^{LM \times LM}$$

Gaussian CRF (GCRF)

The energy function has different nature

GCRF score function

$$F(s \mid X, \theta) = -\frac{1}{2} s^T (W_\theta(X) + \lambda I) s + U_\theta(x) s$$

$$s = [s_{1,1}, \dots, s_{1,L} \mid \dots \mid s_{M,1}, \dots, s_{M,L}]$$

$$U_\theta(X) = [u_\theta^1(X) \mid \dots \mid u_\theta^M(X)] \quad W_\theta(X) \in \mathbb{R}^{LM \times LM}$$

Optimal scores and marginals

$$s^* = (W_\theta(x) + \lambda I)^{-1} U_\theta(X) \quad \text{done efficiently via CG}$$

$$p_j(y_j \mid X, \theta) = \text{softmax}_{y_j}(s_{:,j}^*)$$

Gaussian CRF (GCRF)

The energy function has different nature

GCRF score function

$$F(s \mid X, \theta) = -\frac{1}{2} s^T (W_\theta(X) + \lambda I) s + U_\theta(x) s$$

Prediction

$$\hat{Y} := [\arg \max_k s_{k,1}^*, \dots, \arg \max_k s_{k,L}^*]$$

$$s = [s_{1,1}, \dots, s_{1,L} \mid \dots \mid s_{M,1}, \dots, s_{M,L}]$$

$$U_\theta(X) = [u_\theta^1(X) \mid \dots \mid u_\theta^M(X)] \quad W_\theta(X) \in \mathbb{R}^{LM \times LM}$$

Optimal scores and marginals

$$s^* = (W_\theta(x) + \lambda I)^{-1} U_\theta(X) \quad \text{done efficiently via CG}$$

$$p_j(y_j \mid X, \theta) = \text{softmax}_{y_j}(s_{:,j}^*)$$

Gaussian CRF (GCRF)

The energy function has different nature

GCRF score function

$$F(s \mid X, \theta) = -\frac{1}{2} s^T (W_\theta(X) + \lambda I) s + U_\theta(x) s$$

$$s = [s_{1,1}, \dots, s_{1,L} \mid \dots \mid s_{M,1}, \dots, s_{M,L}]$$

$$U_\theta(X) = [u_\theta^1(X) \mid \dots \mid u_\theta^M(X)] \quad W_\theta(X) \in \mathbb{R}^{LM \times LM}$$

Optimal scores and marginals

$$s^* = (W_\theta(x) + \lambda I)^{-1} U_\theta(X) \quad \text{done efficiently via CG}$$

$$p_j(y_j \mid X, \theta) = \text{softmax}_{y_j}(s_{:,j}^*)$$

Prediction

$$\hat{Y} := [\arg \max_k s_{k,1}^*, \dots, \arg \max_k s_{k,L}^*]$$

Backward pass

$$\frac{\partial \mathcal{L}}{\partial U} = (W + \lambda I)^{-1} \frac{\partial \mathcal{L}}{\partial s^*}, \quad \frac{\partial \mathcal{L}}{\partial \text{vec}(W)} = -\frac{\partial \mathcal{L}}{\partial U} \otimes s^*$$

Gaussian CRF (GCRF)

The energy function has different nature

GCRF score function

$$F(s \mid X, \theta) = -\frac{1}{2} s^T (W_\theta(X) + \lambda I) s + U_\theta(x) s$$

$$s = [s_{1,1}, \dots, s_{1,L} \mid \dots \mid s_{M,1}, \dots, s_{M,L}]$$

$$U_\theta(X) = [u_\theta^1(X) \mid \dots \mid u_\theta^M(X)] \quad W_\theta(X) \in \mathbb{R}^{LM \times LM}$$

Optimal scores and marginals

$$s^* = (W_\theta(x) + \lambda I)^{-1} U_\theta(X) \quad \text{done efficiently via CG}$$

$$p_j(y_j \mid X, \theta) = \text{softmax}_{y_j}(s_{:,j}^*)$$

Prediction

$$\hat{Y} := [\arg \max_k s_{k,1}^*, \dots, \arg \max_k s_{k,L}^*]$$

Backward pass

$$\frac{\partial \mathcal{L}}{\partial U} = (W + \lambda I)^{-1} \frac{\partial \mathcal{L}}{\partial s^*}, \quad \frac{\partial \mathcal{L}}{\partial \text{vec}(W)} = -\frac{\partial \mathcal{L}}{\partial U} \otimes s^*$$

Pairwise potentials

Class-agnostic Potts model

\mathcal{A}_i - pixel embedding

$$W = \begin{bmatrix} 0 & \mathcal{A}\mathcal{A}^T \\ \mathcal{A}\mathcal{A}^T & 0 \end{bmatrix}$$

Online Scaling

Idea

Modify the score function:

$$F_\alpha(Y \mid X, \theta) = F(Y \mid \alpha U_\theta(x), W_\theta(x))$$

Online Scaling

Idea

Modify the score function:

$$F_\alpha(Y \mid X, \theta) = F(Y \mid \alpha U_\theta(x), W_\theta(x))$$

Update scaling factor via grid-search on:

$$\mathcal{L}_\alpha(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(X,Y) \in \mathcal{D}} \mathcal{L}(Y, \alpha U_\theta(X), W_\theta(X))$$

Online Scaling

Idea

Modify the score function:

$$F_\alpha(Y \mid X, \theta) = F(Y \mid \alpha U_\theta(x), W_\theta(x))$$

Update scaling factor via grid-search on:

$$\mathcal{L}_\alpha(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(X,Y) \in \mathcal{D}} \mathcal{L}(Y, \alpha U_\theta(X), W_\theta(X))$$

Algorithm 1 Online Scaling

```
1: initialize  $\theta$ 
2: epoch = 0,  $\alpha = 1$ 
3: for epoch < n_epoch do
4:   for  $(X, Y)$  in training set do
5:     compute  $(U_\theta(X), W_\theta(X))$ 
6:     compute  $\mathcal{L}(Y, \alpha U_\theta(X), W_\theta(X))$ 
7:     update  $\theta$  via an SGD step
8:   end for
9:   epoch += 1
10:   $\mathcal{D}$  := training set or its subset
11:  pick the best  $\alpha_{\text{best}}$  minimizing  $\mathcal{L}_\alpha(\mathcal{D})$  on an  $\alpha$ -grid
12:   $\alpha = \alpha_{\text{best}}$ 
13: end for
```

Online Scaling

Idea

Modify the score function:

$$F_\alpha(Y \mid X, \theta) = F(Y \mid \alpha U_\theta(x), W_\theta(x))$$

Update scaling factor via grid-search on:

$$\mathcal{L}_\alpha(\mathcal{D}) = \frac{1}{|\mathcal{D}|} \sum_{(X,Y) \in \mathcal{D}} \mathcal{L}(Y, \alpha U_\theta(X), W_\theta(X))$$

Remark:

constant scaling \Leftrightarrow init + lr for last layer

Algorithm 1 Online Scaling

```
1: initialize  $\theta$ 
2: epoch = 0,  $\alpha = 1$ 
3: for epoch < n_epoch do
4:   for  $(X, Y)$  in training set do
5:     compute  $(U_\theta(X), W_\theta(X))$ 
6:     compute  $\mathcal{L}(Y, \alpha U_\theta(X), W_\theta(X))$ 
7:     update  $\theta$  via an SGD step
8:   end for
9:   epoch += 1
10:   $\mathcal{D}$  := training set or its subset
11:  pick the best  $\alpha_{\text{best}}$  minimizing  $\mathcal{L}_\alpha(\mathcal{D})$  on an  $\alpha$ -grid
12:   $\alpha = \alpha_{\text{best}}$ 
13: end for
```

Offline Scaling

Explicit relative scaling

$$\hat{U}_\theta(X) = \frac{U_\theta(X)}{\|U_\theta(X)\|} \cdot \alpha, \quad \hat{W}_\theta(X) = \frac{W_\theta(X)}{\|W_\theta(X)\|}$$

Offline Scaling

Explicit relative scaling

$$\hat{U}_\theta(X) = \frac{U_\theta(X)}{\|U_\theta(X)\|} \cdot \alpha, \quad \hat{W}_\theta(X) = \frac{W_\theta(X)}{\|W_\theta(X)\|}$$

or implicit regularizer version

$$\mathcal{R} = \lambda \left(\frac{\|U_\theta(X)\|}{\|W_\theta(X)\|} - \alpha \right)^2$$

Tasks and architectures

Optical Character Recognition

 → command

Unary: LeNet-5 like NN

Pairwise: real-valued matrix of size 26

Tasks and architectures

Optical Character Recognition

 → command

Unary: LeNet-5 like NN

Pairwise: real-valued matrix of size 26

Text Chunking

[NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only # 1.8 billion] [PP in] [NP September].

Figure 6. Chunking example for: He reckons the current account deficit will narrow to only #1.8 billion in September.

Unary: char-level BiLSTM + SENNA + POS + FC

Pairwise: real-value matrix of size N_{tags}

Tasks and architectures

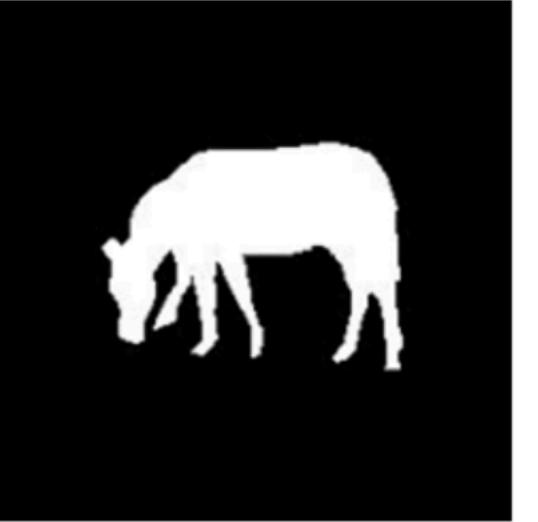
Optical Character Recognition

 command

Unary: LeNet-5 like NN

Pairwise: real-valued matrix of size 26

Binary Segmentation



Text Chunking

Unary and Pairwise: Unet-like FCNNs

[NP He] [VP reckons] [NP the current account deficit] [VP will narrow] [PP to] [NP only # 1.8 billion] [PP in] [NP September].

Figure 6. Chunking example for: He reckons the current account deficit will narrow to only #1.8 billion in September.

Unary: char-level BiLSTM + SENNA + POS + FC

Pairwise: real-value matrix of size N_{tags}

Results

Task	Objective	Stage	Joint	Online	Offline	Offline (Reg)
OCR	Cross-entropy	97.18 ± 0.12	96.48 ± 0.23	97.25 ± 0.11	97.20 ± 0.12	97.14 ± 0.18
	Structured SVM	96.97 ± 0.27	96.43 ± 0.39	97.01 ± 0.24	96.99 ± 0.29	96.91 ± 0.35
	Log-likelihood	97.15 ± 0.13	96.46 ± 0.25	97.17 ± 0.11	97.14 ± 0.13	97.08 ± 0.17
	Cross-entropy (MF)	97.06 ± 0.09	96.62 ± 0.18	97.11 ± 0.08	97.10 ± 0.11	97.01 ± 0.14
Chunking	Cross-entropy	89.52 ± 0.29	87.92 ± 0.38	89.56 ± 0.26	89.53 ± 0.28	89.45 ± 0.34
	Structured SVM	89.21 ± 0.49	87.64 ± 0.57	89.25 ± 0.44	89.19 ± 0.51	89.13 ± 0.55
	Log-likelihood	89.43 ± 0.31	87.85 ± 0.37	89.48 ± 0.27	89.45 ± 0.29	89.35 ± 0.36
	Cross-entropy (MF)	89.25 ± 0.26	87.68 ± 0.29	89.34 ± 0.22	89.27 ± 0.30	89.15 ± 0.28
Bin. segm.	Cross-entropy	86.51 ± 0.19	85.62 ± 0.31	86.56 ± 0.17	86.60 ± 0.21	86.48 ± 0.25

Conclusion

- Conjecture: scaling matters
- We propose online and offline scaling algorithms
- And demonstrate their efficiency on three different tasks