Distributed Training Strategies for the Structured Perceptron

Ryan McDonald Keith Hall Gideon Mann Google, Inc., New York / Zurich

Presenter: Mahsa Yarmohammadi

- An algorithm for supervised classification
- Classify one input into one of several outputs
- Example in POS tagging:

input:

"fruit flies fast" possible outputs: *"NN NNS VB" "NN VB RB" "NN NNS JJ"*

- Supervised: input-output training instances $(\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{T}$
- For all training instances:
 - Predict the best output (y')
 - Penalize features of \mathbf{y}'
 - Reward features of the truth (y_t)

Predict the best output ('y')

$$\mathbf{y}' = \operatorname{arg\,max}_{\mathbf{y}'} \mathbf{w}^{(k)} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$$

Predict the best output (y')



• Predict the best output ('y')

$$\mathbf{y}' = \operatorname{arg\,max}_{\mathbf{y}'} \mathbf{w}^{(k)} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$$

• Reward/Penalize features

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$$

Perceptron(
$$\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$$
)
 $\mathbf{w}^{(0)} = \mathbf{0}; \ k = 0$
for $n : 1..N$
for $t : 1..T$
Let $\mathbf{y}' = \arg \max_{\mathbf{y}'} \mathbf{w}^{(k)} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$
if $\mathbf{y}' \neq \mathbf{y}_t$
 $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \mathbf{y}')$
 $k = k + 1$
return $\mathbf{w}^{(k)}$

- Can be costly to train
 - increasing size of training data
 - complexity is proportional to inference
 - non-linear increasing complexity with sequence length
- Distributed training strategies to reduce training time

Distributed Perceptron

- Two strategies:
 - Parameter Mixing
 - Iterative Parameter Mixing
- Evaluate on two problems:
 - Named Entity Recognition (NER)
 - Dependency Parsing

1) Divide training data \mathcal{T} into S shards

2) Train perceptron on each shard in parallel

3) Final \mathbf{w} is a weighted mixture of perceptrons

1) Divide training data ${\cal T}~{\rm into}S~{\rm shards}$

2) Train perceptron on each shard in parallel mapper
3) Final w is a weighted mixture of perceptrons
reducer

PerceptronParamMix(
$$\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$$
)

Shard \mathcal{T} into S pieces $\mathcal{T} = \{\mathcal{T}_1, \dots, \mathcal{T}_S\}$ $\mathbf{w}^{(i)} = \text{Perceptron}(\mathcal{T}_i)$ $\mathbf{w} = \sum_i \mu_i \mathbf{w}^{(i)}$ return \mathbf{w}

- Advantages:
 - straight-forward
 - scalable to extremely large data sets
 - resource (network usage) efficient
- Disadvantage:
 - sub-optimal: does not necessarily return a separating weight vector

Separable Training Set

• \mathcal{T} is separable with margin γ if there is $\mathbf{u}(||\mathbf{u}|| = 1)$

$$\mathbf{u} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{u} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}') \geq \gamma$$

for all $(\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{T}$, and for all $\mathbf{y}' \in \mathcal{Y}_t$ such that $\mathbf{y}' \neq \mathbf{y}_t$

Separable Training Set

• \mathcal{T} is separable with margin γ if there is $\mathbf{u}(||\mathbf{u}|| = 1)$

$$\mathbf{u} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{u} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}') \geq \gamma$$

for all $(\mathbf{x}_t, \mathbf{y}_t) \in \mathcal{T}$, and for all $\mathbf{y}' \in \mathcal{Y}_t$ such that $\mathbf{y}' \neq \mathbf{y}_t$

- If ${\mathcal T}$ is separable, then the perceptron
- 1) will converge in a finite time
- 2) will produce a separating ${\bf w}$

1) Divide training data \mathcal{T} into S shards

2) Train single-epoch perceptron on each shard

3) Mix the weights

4) Re-send the new weight w to each shard, go to 2

1) Divide training data ${\cal T}~{\rm into}S~{\rm shards}$

2) Train single-epoch perceptron on each shard mapper
3) Mix the weights reducer
4) Re-send the new weight w to each shard, go to 2

PerceptronIterParamMix($\mathcal{T} = \{(\mathbf{x}_t, \mathbf{y}_t)\}_{t=1}^{|\mathcal{T}|}$) Shard \mathcal{T} into S pieces $\mathcal{T} = \{\mathcal{T}_1, \dots, \mathcal{T}_S\}$ $\mathbf{w} = \mathbf{0}$ for n : 1..N $\mathbf{w}^{(i,n)} = \text{OneEpochPerceptron}(\mathcal{T}_i, \mathbf{w})$ $\mathbf{w} = \sum_i \mu_{i,n} \mathbf{w}^{(i,n)}$ return \mathbf{w}

```
OneEpochPerceptron(\mathcal{T}, \mathbf{w}^*)

\mathbf{w}^{(0)} = \mathbf{w}^*; k = 0

for t : 1..T

Let \mathbf{y}' = \arg \max_{\mathbf{y}'} \mathbf{w}^{(k)} \cdot \mathbf{f}(\mathbf{x}_t, \mathbf{y}')

if \mathbf{y}' \neq \mathbf{y}_t

\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \mathbf{y}')

k = k + 1

return \mathbf{w}^{(k)}
```

- Advantage:
 - optimal: returns a *separating* weight vector
 - significantly reduces training time
- Disadvantage:
 - increasing network usage

Experiments

- Standard/Averaged perceptron in 4 systems:
 - Serial All Data
 - Serial Sub Sampling
 - Parallel (Parameter Mixing)
 - Parallel (Iterative Parameter Mixing)

Experiments

- Standard/Averaged perceptron in 4 systems:
 - Serial All Data
 - Serial Sub Sampling
 - Parallel (PM)
 - Parallel (IPM)
- Two Problems:
 - NER CoNLL 2033 shared task
 - Dependency parsing Prague Dependency Treebank





• Training on a single shard performs worse than all data



PM performs better than single-shard but worse than all data



IPM performs as good as/better than all data



 PM and IPM return better classifiers much quicker than all data

NER – Averaged perceptron



The same patterns hold for averaged perceptron



	Reg. Perceptron	Avg. Perceptron
	F-measure	F-measure
Serial (All Data)	85.8	88.2
Serial (Sub Sampling)	75.3	76.6
Parallel (Parameter Mix)	81.5	81.6
Parallel (Iterative Parameter Mix)	87.9	88.1

Dep. Parsing – Standard perceptron



Dep. Parsing – Averaged perceptron





	Reg. Perceptron	Avg. Perceptron
	Unlabeled Attachment Score	Unlabeled Attachment Score
Serial (All Data)	81.3	84.7
Serial (Sub Sampling)	77.2	80.1
Parallel (Iterative Parameter Mix)	83.5	84.5



- # shards
 - 10
 - 100
- Mixing strategy
 - uniform ($\mu_{i,n}=1/S$)
 - error based ($\mu_{i,n} = k_{i,n}/k_n$)
- # epochs (N)
 - -0..50





• Mixing strategy makes little difference



Increasing # shards slows down convergence when # epochs increases



 Trade-off between slower convergence and quicker ephocs

Conclusions

- Iterative parameter mixing is
 - guaranteed to separate the data
 - significantly reduces training time
- Trade-off between distributing computation and slower convergence relative to # shards
- Theoretical proofs of convergence in the paper