Training Strategies

CS 6355: Structured Prediction



So far we saw

- What is structured output prediction?
- Different ways for modeling structured prediction
 Conditional random fields, factor graphs, constraints
- What we only occasionally touched upon:
 - Algorithms for training and inference
 - Viterbi (inference in sequences)
 - Structured perceptron (training in general)

Rest of the semester

• Strategies for training

- Structural SVM
- Stochastic gradient descent
- More on local vs. global training
- Algorithms for inference
 - Exact inference
 - "Approximate" inference
 - Formulating inference problems in general
- Latent/hidden variables, representations and such

Up next

- Structural Support Vector Machine
 - How it naturally extends multiclass SVM
- Empirical Risk Minimization
 - Or: how structural SVM and CRF are solving very similar problems
- Training Structural SVM via stochastic gradient descent
 - And some tricks

Where are we?

- Structural Support Vector Machine
 - How it naturally extends multiclass SVM
- Empirical Risk Minimization
 - Or: how structural SVM and CRF are solving very similar problems
- Training Structural SVM via stochastic gradient descent
 - And some tricks

Recall: Binary and Multiclass SVM

- Binary SVM
 - Maximize margin
 - Equivalently,

Minimize norm of weights such that the closest points to the hyperplane have a score §1

- Multiclass SVM
 - Each label has a different weight vector (like one-vs-all)
 - Maximize multiclass margin
 - Equivalently,

Minimize total norm of the weights such that the true label is scored at least 1 more than the second best one

Multiclass SVM in the separable case

Recall hard binary SVM

We have a data set $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

 $\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$ s.t. $\forall i, \quad y_i \mathbf{w}^T \mathbf{x}_i \ge 1$

Multiclass SVM in the separable case



Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

Features sum over the parts

Modeling





Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

Features sum over the parts

Modeling

Data





Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

• Features sum over the parts

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{p \in \text{parts}(\mathbf{x})} \Phi_p(\mathbf{x}, \mathbf{y}_p)$$

We also have a data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

What we want from training (following the multiclass idea)

For each training example $(\mathbf{x}_i, \mathbf{y}_i)$:

Modeling

Data

- The annotated structure \mathbf{y}_i gets the highest score among all structures

Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

• Features sum over the parts

Modeling

Data

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{p \in \text{parts}(\mathbf{x})} \Phi_p(\mathbf{x}, \mathbf{y}_p)$$

We also have a data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

- The annotated structure \mathbf{y}_i gets the highest score among all structures
- Or to be safe, \mathbf{y}_i gets a score that is at least one more than all other structures

$$\forall \mathbf{y} \neq \mathbf{y}_i, \qquad \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}) + 1$$

Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

• Features sum over the parts

Modeling

Data

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{p \in \text{parts}(\mathbf{x})} \Phi_p(\mathbf{x}, \mathbf{y}_p)$$

We also have a data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

- The annotated structure \mathbf{y}_i gets the highest score among all structures
- Or to be safe, \mathbf{y}_i gets a score that is at least one more than all other structures

$$\forall \mathbf{y} \neq \mathbf{y}_i, \qquad \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}) + 1$$

Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

• Features sum over the parts

Modeling

Data

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{p \in \text{parts}(\mathbf{x})} \Phi_p(\mathbf{x}, \mathbf{y}_p)$$

We also have a data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

- The annotated structure \mathbf{y}_i gets the highest score among all structures
- Or to be safe, \mathbf{y}_i gets a score that is at least one more than all other structures

$$\forall \mathbf{y} \neq \mathbf{y}_i, \qquad \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}) + 1$$

Suppose we have some definition of a structure (a factor graph)

And feature definitions for each factor (i.e. "part") p as $\Phi_p(\mathbf{x}, \mathbf{y}_p)$

Remember: we can talk about the feature vector for the entire structure

• Features sum over the parts

Modeling

Data

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_{p \in \text{parts}(\mathbf{x})} \Phi_p(\mathbf{x}, \mathbf{y}_p)$$

We also have a data set $D = \{(\mathbf{x}_i, \mathbf{y}_i)\}$

- The annotated structure \mathbf{y}_i gets the highest score among all structures
- Or to be safe, \mathbf{y}_i gets a score that is at least one more than all other structures

$$\forall \mathbf{y} \neq \mathbf{y}_i, \qquad \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}) + 1$$

Maximize margin

s.t. Score for gold \geq Score for other +1 For every training example

Maximize margin





Problem?



Problem





Problem





Other structure A: Only one mistake



Other structure B: Fully incorrect



Problem



Structure B has is more wrong, but this formulation will be happy if *both* A & B are scored one less than gold! No partial credit!



Other structure A: Only one mistake



Other structure B: Fully incorrect





Hamming distance between structures: Counts the number of differences between them



because the Hamming distance of **y** and itself is zero ²⁵

$$\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

s.t. $\mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in D, \forall \mathbf{y}$

- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores

 $\min_{\mathbf{w}} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad \longleftarrow \quad \text{Maximize margin by minimizing norm of } \mathbf{w}$

s.t. $\mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in D, \forall \mathbf{y}$

- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



- It is okay for a structure that is close (in Hamming sense) to the true one to get a score that is close to the true structure
- Structures that are very different from the true structure should get much lower scores



Problem?



Problem? What if the data is not separable?

• What if these constraints are not satisfied for any **w** for a given dataset?





 $\min_{\mathbf{w},\xi} \qquad \qquad \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_i \xi_i$

s.t. $\mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \ge \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) - \xi_i \quad \forall (\mathbf{x}_i, \mathbf{y}_i) \in D, \forall \mathbf{y}$



For every labeled example, and every competing structure



For every labeled example, and every competing structure, the score for the ground truth should be greater than the score for the competing structure by the Hamming distance between them



Slack variables allow some examples to be misclassified.



Slack variables allow some examples to be misclassified.



Slack variables allow some examples to be misclassified.

Minimizing the slack forces this to happen as few times as possible



Slack variables allow some examples to be misclassified.

Minimizing the slack forces this to happen as few times as possible

Questions?

Structural SVM



Structural SVM





Comments

- Other slightly different formulations exist
 - Generally same principle
- Multiclass is a special case of structure
 - Structural SVM strictly generalizes multiclass SVM Exercise: Work it out
- Can be seen as minimizing structured version of hinge loss
 - Remember empirical risk minimization?
- Learning as optimization
 - We have framed the optimization problem
 - We haven't seen how it can be solved yet
 - That is, we don't have a learning algorithm yet

Where are we?

- Structural Support Vector Machine
 - How it naturally extends multiclass SVM
- Empirical Risk Minimization
 - Or: how structural SVM and CRF are solving very similar problems
- Training Structural SVM via stochastic gradient descent
 - And some tricks

Broader picture: Learning as loss minimization

- Collect some annotated data. More is generally better
- Pick a hypothesis class (also called model)
 - Decide how the score decomposes over the parts of the output
- Choose a loss function
 - Decide on how to penalize incorrect decisions
- Learning = minimize empirical risk + regularizer
 - Typically an optimization procedure needed here

This must look familiar. We have seen this before for binary classification!

Empirical risk minimization

- Suppose the function *Loss* scores the quality of a prediction with respect to the true structure
 - Loss(f(x), y) tells us how good f is for this x by comparing it against y
- Evaluate the quality of the predictor *f* by averaging over the unknown distribution P that generates data
 - Expected risk: $\mathbb{E}[Loss(f(\mathbf{x}), \mathbf{y})]$
- We don't know P, so use the empirical risk $\frac{1}{N} \sum_{i} Loss(f(\mathbf{x}_i), \mathbf{y}_i)$ possibly with regularizer

Learning: Minimizing regularized risk; various algorithms exist

The loss function zoo: binary classification



The loss function zoo



• Structural SVM $\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i} \max_{\mathbf{y}} \left(\mathbf{w}^T \Phi(\mathbf{x}_i, \mathbf{y}) + \Delta(\mathbf{y}, \mathbf{y}_i) - \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y}_i) \right)$

• Conditional Random Field (via the maximum a posteriori criterion)

$$\min_{\mathbf{w}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i} -\log P(\mathbf{y}_i \mid \mathbf{x}_i, \mathbf{w})$$
Where *P* is defined as

where P is defined as $evn(\mathbf{w}^T \Phi(\mathbf{x}, \mathbf{y}))$

$$P(\mathbf{y}_i \mid x_i, w) = \frac{\exp(\mathbf{w} \cdot \Phi(\mathbf{x}_i, \mathbf{y}_i))}{Z(\mathbf{x}_i, \mathbf{w})}$$











Structured Perceptron

$$\min_{\mathbf{w}} \sum_{i} \max_{\mathbf{y}} \left(\mathbf{w}^{T} \Phi(\mathbf{x}_{i}, \mathbf{y}) - \mathbf{w}^{T} \phi(\mathbf{x}_{i}, \mathbf{y}_{i}) \right)$$



Summary

- Different structured training objectives are really different loss functions
- The structured versions of hinge, log and Perceptron losses all involve inference
 - Hinge, Perceptron: Solve a maximization problem
 - Log: Solve an expectation problem
- Learning as stochastic optimization, even for structures
 - But, computing the loss (and the gradient) can be expensive