General Formulations for Structures: Constrained Conditional Models

CS 6355: Structured Prediction



Where are we?

- Graphical models
 - Bayesian Networks
 - Markov Random Fields
- Formulations of structured output
 - Joint models
 - Markov Logic Network
 - Conditional models
 - Conditional Random Fields (again)
 - Constrained Conditional Models

Outline

- Consistency of outputs and the value of inference
- Constrained conditional models via an example
- Hard constraints and Integer Programs
- Soft constraints

Outline

- Consistency of outputs and the value of inference
- Constrained conditional models via an example
- Hard constraints and Integer Programs
- Soft constraints

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B Suppose, for some example, we have a sequence with three steps



Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

Suppose, for some example, we have a sequence with three steps We can define a linear chain factor graph with unary and binary terms



Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

Suppose, for some example, we have a sequence with three steps We can define a linear chain factor graph with unary and binary terms



This is short hand for this factor graph. Going ahead, we will not show the inputs. We will assume that we have conditional models.

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B



Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output



Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output



Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output

How can we do that? Which decisions interact in this condition?



Possible outputs AAA AAB ABA BAA BAA BAB BBA BBA BBB

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output



Possible outputs AAA ABB ABA BAA BAA BAB BBA BBB

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output



This potential function ensures the condition

y1	y2	у3	f
А	А	А	1
Α	А	В	1
Α	В	А	1
Α	В	В	0
В	А	А	1
В	А	В	0
В	В	А	0
В	В	В	0

Or: How to introduce knowledge into prediction

Suppose we have a sequence labeling problem where the outputs can be one of A or B

We want to add a condition:

There should be no more than one B in the output



But the standard CRF learning does not allow for potential functions to be set manually

Should we learn what we can write down easily?

Especially for such large, computationally cumbersome factors

This potential function ensures the condition

y1	y2	у3	f
А	А	А	1
А	А	В	1
A	В	А	1
А	В	В	0
В	А	А	1
В	А	В	0
В	В	А	0
В	В	В	0

Another look at learning and inference

• Inference: A global decision comprising of multiple local decisions and their inter-dependencies

If there were no inter-dependencies between decisions, we could as well treat these as independent prediction problems

Another look at learning and inference

• Inference: A global decision comprising of multiple local decisions and their inter-dependencies

If there were no inter-dependencies between decisions, we could as well treat these as independent prediction problems

• Does **learning** need to be global too?

Recall: Local vs. global learning

- Global learning: Learn with inference
- Local learning: Learn the local decisions independently and piece them together
 - Maybe we can learn <u>sub-structures</u> independently and piece them together

Typical updates in global learning

Stochastic gradient descent update for CRF

- For a training example $(\mathbf{x}_i, \mathbf{y}_i)$: $\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\Phi(\mathbf{x}_i, \mathbf{y}_i) - E_{\mathbf{y} \sim P(\cdot | \mathbf{x}_i, \mathbf{w})} [\Phi(\mathbf{x}_i, \mathbf{y}_i)])$

Typical updates in global learning

Stochastic gradient descent update for CRF

- For a training example
$$(\mathbf{x}_i, \mathbf{y}_i)$$
:
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\Phi(\mathbf{x}_i, \mathbf{y}_i) - E_{\mathbf{y} \sim P(\cdot | \mathbf{x}_i, \mathbf{w})} [\Phi(\mathbf{x}_i, \mathbf{y}_i)])$

Structured perceptron

- For a training example
$$(\mathbf{x}_i, \mathbf{y}_i)$$
:
 $\hat{y} = \max_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y})$
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \hat{\mathbf{y}}))$

Typical updates in global learning

Stochastic gradient descent update for CRF

- For a training example
$$(\mathbf{x}_i, \mathbf{y}_i)$$
:
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\Phi(\mathbf{x}_i, \mathbf{y}_i) - E_{\mathbf{y} \sim P(\cdot | \mathbf{x}_i, \mathbf{w})} [\Phi(\mathbf{x}_i, \mathbf{y}_i)])$

Structured perceptron

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

$$\hat{y} = \max_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}_i, \mathbf{y})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha_t (\Phi(\mathbf{x}_i, \mathbf{y}_i) - \Phi(\mathbf{x}_i, \hat{\mathbf{y}}))$$

Some form of inference needed within the innermost loop

Why local learning?

- Global learning may not always be feasible
 - Practical concerns

Global learning can be computationally prohibitive because of inference within the training loop

- Data issues
 - What if we don't have a single dataset that is fully annotated with a structure
 - Instead, we have multiple datasets, each with a subset of the structures
- But inference *at deployment* can still be global
 - Recall the discussion of local vs global learning from multiclass classification

In the one-vs-all case, learning was local, but prediction was not.

Combining local classifiers

Inference can "glue" together local decisions

And enforce *global* coherence



Setting

Output: Nodes and edges are labeled and the blue and orange edges form a tree

Goal: Find the highest scoring tree

Combining local classifiers

Inference can "glue" together local decisions

- And enforce *global* coherence



Setting

Output: Nodes and edges are labeled and the blue and orange edges form a tree

Goal: Find the highest scoring tree



We need to ensure that the colored edges form a valid output (i.e. a tree)



Combining local classifiers

Inference can "glue" together local decisions

- And enforce *global* coherence



Setting

Output: Nodes and edges are labeled and the blue and orange edges form a tree

Goal: Find the highest scoring tree



We need to ensure that the colored edges form a valid output (i.e. a tree)



Combining local classifiers

Inference can "glue" together local decisions

- And enforce *global* coherence



Setting

Output: Nodes and edges are labeled and the blue and orange edges form a tree

Goal: Find the highest scoring tree



We need to ensure that the colored edges form a valid output (i.e. a tree)



Constraints at prediction time

Introduce additional information

Might not have been available at training time

Add domain knowledge

Examples:

- "All part-of-speech tag sequences should contain a verb"
- "Every bicycle should have at least one wheel"
- "In any window of size six, at least one of the labels should be a B"

Enforce coherence into the set of local decisions

Examples:

- "the collection of decisions should form a tree"
- "the collection of parts recognized should form a valid bicycle"

Outline

- Consistency of outputs and the value of inference
- Constrained conditional models via an example
- Hard constraints and Integer Programs
- Soft constraints

Constrained Conditional Model

Inference consists of two components

- 1. Local classifiers
 - <u>Important</u>: These may be a collection of structures themselves
 - These are trained models
- 2. A set of constraints that restrict the space of joint assignments of the local classifiers
 - Background knowledge





Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i} \mathbf{w}^{T} \phi_{i}(\mathbf{x}, \mathbf{y}_{i})$$



Or equivalently,





Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i} \mathbf{w}^{T} \phi_{i}(\mathbf{x}, \mathbf{y}_{i})$$

Or equivalently (here),

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \mathbf{w}^{T} \phi_{E}(\mathbf{x}, \mathbf{y}_{i}) + \sum_{i=1}^{2} \mathbf{w}^{T} \phi(\mathbf{y}_{i}, \mathbf{y}_{i+1})$$



$$\operatorname{argmax}_{\mathbf{y}} \sum_{i=1}^{3} \mathbf{w}^{T} \phi_{E}(\mathbf{x}, \mathbf{y}_{i}) + \sum_{i=1}^{2} \mathbf{w}^{T} \phi(\mathbf{y}_{i}, \mathbf{y}_{i+1})$$

$$\downarrow$$
Emissions
$$\downarrow$$
Transitions
















Prediction in a CCM Suppose the outputs can be one of A or B Typically, we have $\underset{y}{\operatorname{argmax}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_1)$ Or equivalently, $\frac{3}{2}$

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{N} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{N} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

One group of indicators per factor One score per indicator



One group of indicators per factor One score per indicator This expression explicitly enumerates *every* decision that we need to make to build the final output











Prediction in a CCM
Suppose the outputs can be one of A or B
Typically, we have
$$\operatorname{argmax} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$$

y
 $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_1)$
 $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_2)$
 $\mathbf{w}^$

Н

Some decisions can not exist together

Prediction in a CCM
Suppose the outputs can be one of A or B
Typically, we have
$$\operatorname{argmax} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$$

 $\mathbf{y}_{1=A}$ score($y_1=A$) + $l_{y_{1=B}}$ score($y_1=B$)
 $l_{y_{2=A}}$ core($y_2=A$) + $l_{y_{2=B}}$ score($y_2=B$)
+ $l_{y_{1=B} \text{ and } y_{2=A}}$ score($y_1=A, y_2=A$) + $l_{y_{1=A} \text{ and } y_{2=B}}$ score($y_1=A, y_2=B$)
+ $l_{y_{1=B} \text{ and } y_{2=A}}$ score($y_1=B, y_2=A$) + $l_{y_{1=B} \text{ and } y_{2=B}}$ score($y_1=B, y_2=B$)
+ $l_{y_{2=A} \text{ and } y_{3=A}}$ score($y_2=A, y_3=A$) + $l_{y_{2=B} \text{ and } y_{3=B}}$ score($y_2=B, y_3=A$) + $l_{y_{2=B} \text{ and } y_{3=B}}$ score($y_2=B, y_3=A$) + $l_{y_{2=B} \text{ and } y_{3=B}}$ score($y_2=B, y_3=B$)
Some decisions can not exist together
If $l_{y_2=A}$ = 1 then $l_{y_1=B}$ and $y_{3=A}$ = 0
If $l_{y_2=A}$ = 1 then $l_{y_1=B}$ and $y_{2=B}$ = 0



$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{5} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

But we should only consider valid label sequences

Prediction in a CCM Suppose the outputs can be one of A or B y $w^{T}\phi_{T}(y_{1}, y_{2})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{E}(x, y_{1})$ $w^{T}\phi_{E}(x, y_{2})$ $w^{T}\phi_{E}(x, y_{2})$ $w^{T}\phi_{E}(x, y_{3})$

Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

Predict with constraints

• Each y_i can either be a A or a B

Prediction in a CCM Suppose the outputs can be one of A or B y $w^{T}\phi_{T}(y_{1}, y_{2})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{E}(x, y_{1})$ $w^{T}\phi_{E}(x, y_{2})$ $w^{T}\phi_{E}(x, y_{2})$ $w^{T}\phi_{E}(x, y_{3})$

Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

Predict with constraints

• Each y_i can either be a A or a B: $\forall i, I_{y_i=A} + I_{y_i=B} = 1$

Prediction in a CCM Suppose the outputs can be one of A or B Typically, we have $\underset{y}{\operatorname{argmax}} \mathbf{w}^{T} \phi(\mathbf{x}, \mathbf{y})$ $\mathbf{w}^{T} \phi_{\mathsf{E}}(\mathbf{x}, y_{1})$ $\mathbf{w}^{T} \phi_{\mathsf{E}}(\mathbf{x}, y_{2})$ $\mathbf{w}^{T} \phi_{\mathsf{E}}(\mathbf{x}, y_{2})$ $\mathbf{w}^{T} \phi_{\mathsf{E}}(\mathbf{x}, y_{2})$ $\mathbf{w}^{T} \phi_{\mathsf{E}}(\mathbf{x}, y_{2})$

Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

Predict with constraints

- Each y_i can either be a A or a B: $\forall i, I_{y_i=A} + I_{y_i=B} = 1$
- The emission decisions and the transition decisions should agree

We can write this using linear constraints over the indicator variables

Prediction in a CCM Suppose the outputs can be one of A or B Typically, we have $\underset{y}{\operatorname{argmax}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_1)$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_2)$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_2)$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_2)$ $\mathbf{w}^T \phi_{\mathsf{E}}(\mathbf{x}, y_2)$

Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

Predict with constraints

- Each y_i can either be a A or a B: $\forall i, I_{y_i=A} + I_{y_i=B} = 1$
- The emission decisions and the transition decisions should agree
- There should be no more than one B in the output

We could add extra knowledge that was not present at training time

Prediction in a CCM Suppose the outputs can be one of A or B y $w^{T}\phi_{T}(y_{1}, y_{2})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{T}(y_{2}, y_{3})$ $w^{T}\phi_{E}(\mathbf{x}, y_{1})$ $w^{T}\phi_{E}(\mathbf{x}, y_{2})$ $w^{T}\phi_{E}(\mathbf{x}, y_{2})$ $w^{T}\phi_{E}(\mathbf{x}, y_{3})$

Or equivalently,

$$\underset{\mathbf{y}}{\operatorname{argmax}} \sum_{i=1}^{3} \sum_{l \in \{A,B\}} I_{y_{i}=l} \cdot \mathbf{w}^{T} \phi_{E}(\mathbf{x}, l) + \sum_{i=1}^{2} \sum_{l_{1}, l_{2} \in \{A,B\}} I_{y_{i}=l_{1} \wedge y_{i+1}=l_{2}} \mathbf{w}^{T} \phi(l_{1}, l_{2})$$

Predict with constraints

- Each y_i can either be a A or a B: $\forall i, I_{y_i=A} + I_{y_i=B} = 1$
- The emission decisions and the transition decisions should agree
- There should be no more than one B in the output $I_{y_1=B} + I_{y_2=B} + I_{y_3=B} \le 1$

Questions?

with hard constraints

- Global joint inference
 - Treat the output as a collection of decisions (one per factor/part)
 - Each decision associated with a score
 - In addition, allow arbitrary constraints involving the decisions
- Constraints
 - Inject domain knowledge into the prediction
 - Can be stated as logical statements
 - Can be transformed into linear inequalities in terms of the decision variables
- No comment about learning
 - Learn the scoring functions locally or globally

Outline

- Consistency of outputs and the value of inference
- Constrained conditional models via an example
- Hard constraints and Integer Programs
- Soft constraints

$$\underset{\mathbf{y}}{\operatorname{arg\,max}} \sum_{p \in \operatorname{Parts}(\mathbf{x})} \sum_{l \in \operatorname{Labels}} I_{y_p = l} \cdot \operatorname{score}(\mathbf{x}, l)$$

Such that **y** is feasible





- Can be written as linear (in)equalities in the I's
- I's can only be 0 or 1

This is an Integer Linear Program



- Can be written as linear (in)equalities in the I's
- I's can only be 0 or 1

This is an Integer Linear Program

MAP inference can be written as integer linear programs (ILPs)

У



- Can be written as linear (in)equalities in the I's
- I's can only be 0 or 1

This is an Integer Linear Program

MAP inference can be written as integer linear programs (ILPs)

But solving ILPs is NP-complete in the worst case

- Use an off-the-shelf solver (Gurobi) and hope for the best
- **Or not**: Write a specialized search algorithm if we know more about the problem (exact or approximate)
- We will see examples of these

Questions?

ILP inference

- Can introduce domain knowledge in the form of constraints
 - Any Boolean expression over the inference variables can be written as linear inequalities
- A uniform language for formulating and reasoning about inference
- Have not made the problem any easier to solve
 - By allowing easy addition of constraints, it may be simple to write down provably intractable inference formulations
 - (Off-the-shelf solvers seem to work admirably!)

Outline

- Consistency of outputs and the value of inference
- Constrained conditional models via an example
- Hard constraints and Integer Programs
- Soft constraints

Constraints may not always hold

"Every car should have a wheel"



General case: with soft constraints



The model score for the structure

General case: with soft constraints

$$\arg \max_{\mathbf{y}} \mathbf{w}^T \phi(\mathbf{x}, \mathbf{y})$$

The model score for the structure

Suppose we have a collection of "soft constraints" C_1, C_2, \cdots each associated with penalties for violation ρ_1, ρ_2, \cdots

That is: A constraint C_k is a Boolean expression over the output. An assignment that violates this constraint has to pay a penalty of ρ_k

General case: with soft constraints



General case: with soft constraints



Constrained conditional models: review

- Write down conditions that the output need to satisfy
 - Constraints are effectively factors in a factor graph whose potential functions are fixed
- Different learning regimes
 - Train with the constraints or without
 - Remember: constraint penalties are fixed in either case
- Prediction
 - Can write the inference formulation as an integer linear program
 - Can solve it with an off-the-shelf solver (or not!)
- Extension
 - Soft constraints: Constraints that don't always need to hold

Structured prediction: General formulations

- Graphical models
 - Bayesian Networks
 - Markov Random Fields
- Formulations of structured output
 - Joint models
 - Markov Logic Network
 - Conditional models
 - Conditional Random Fields (again)
 - Constrained Conditional Models

All the discussion so far has been about representation.

Forthcoming lectures: Algorithms for training and inference