# *From* Structured Prediction *to* Inverse Reinforcement Learning

Hal Daumé III

Computer Science
University of Maryland

me@hal3.name

A Tutorial at ACL
Uppsala Universitet

Sunday, 11 July 2010

# NLP as transduction

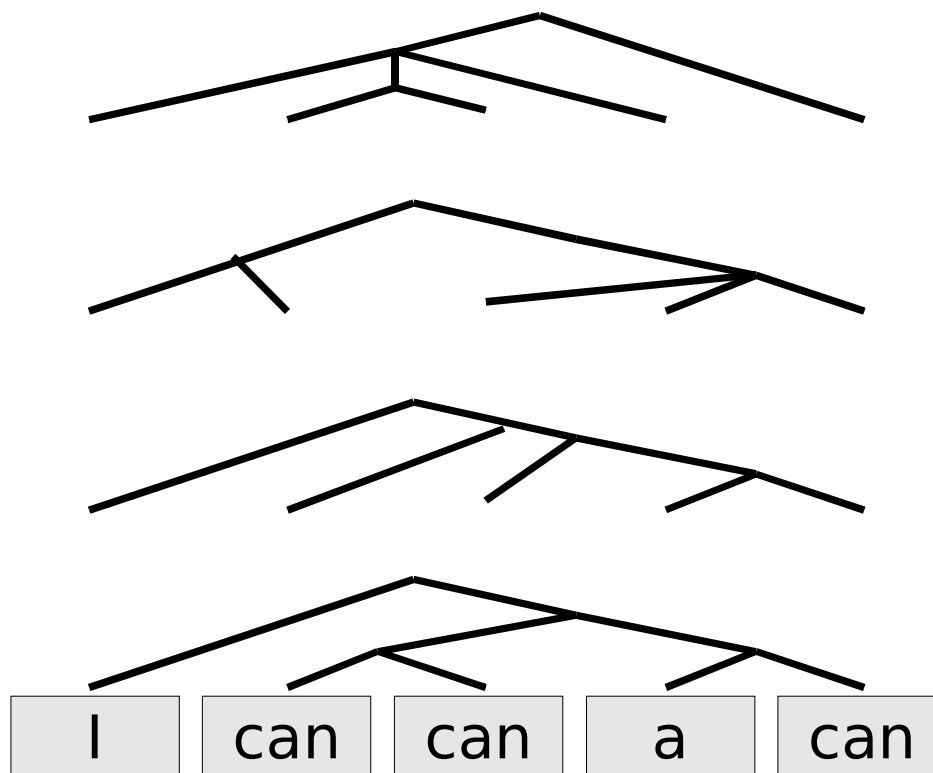| Task | Input | Output |
|---|---|---|
| Machine Translation | Ces deux principes se tiennent à la croisée de la philosophie, de la politique, de l'économie, de la sociologie et du droit. | Both principles lie at the crossroads of philosophy, politics, economics, sociology, and law. |
| Document Summarization | Argentina was still obsessed with the Falkland Islands even in 1994, 12 years after its defeat in the 74-day war with Britain. The country's overriding foreign policy aim continued to be winning sovereignty over the islands. | The Falkland islands war, in 1982, was fought between Britain and Argentina. |
| Syntactic Analysis | The man ate a big sandwich. | The man ate a big sandwich |
| ...many more... | | |

# Structured prediction 101

Learn a function mapping inputs to complex outputs:

$$f : X \rightarrow Y$$

Input Space    Decoding    Output Space



| I | can | can | a | can |

Hal Daumé III (me@hal3.name)    SP2IRL @ ACL2010

# Why is structure important?

➢ Correlations among outputs
  ➢ Determiners often precede nouns
  ➢ Sentences usually have verbs

➢ Global coherence
  ➢ It just *doesn't make sense* to have three determiners next to each other

➢ My objective (aka "loss function") forces it
  ➢ Translations should have good sequences of words
  ➢ Summaries should be coherent

# Outline: Part I

- What is Structured Prediction?
- Refresher on Binary Classification
  - What does it mean to learn?
  - Linear models for classification
  - Batch versus stochastic optimization
- From Perceptron to Structured Perceptron
  - Linear models for Structured Prediction
  - The "argmax" problem
  - From Perceptron to margins
- Learning to Search
  - Incremental Parsing
  - Search-based Structured Prediction

# Outline: Part II

- Refresher on Reinforcement Learning
  - Markov Decision Processes
  - Q learning
- Apprenticeship Learning
  - Inverse RL
  - Apprenticeship Learning via IRL
- Inverse Optimal Control and A* Search
  - Maximum Margin Planning
  - Learning to Search

- Discussion

# Refresher on Binary Classification

# What does it mean to learn?

- Informally:
  - to predict the future based on the past

- Slightly-less-informally:
  - to take *labeled examples* and construct a function that will label them as a human would

- Formally:
  - Given:
    - A fixed unknown distribution D over X*Y
    - A loss function over Y*Y
    - A finite sample of (x,y) pairs drawn i.i.d. from D
  - Construct a function f that has low expected loss with respect to D

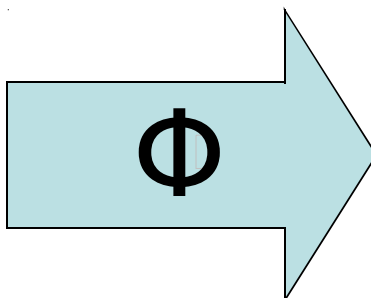# Feature extractors

➢ A feature extractor Φ maps examples to vectors

```
Dear Sir.

First, I must solicit
your confidence in
this transaction,
this is by virture of
its nature as being
utterly confidencial
and top secret. …
```
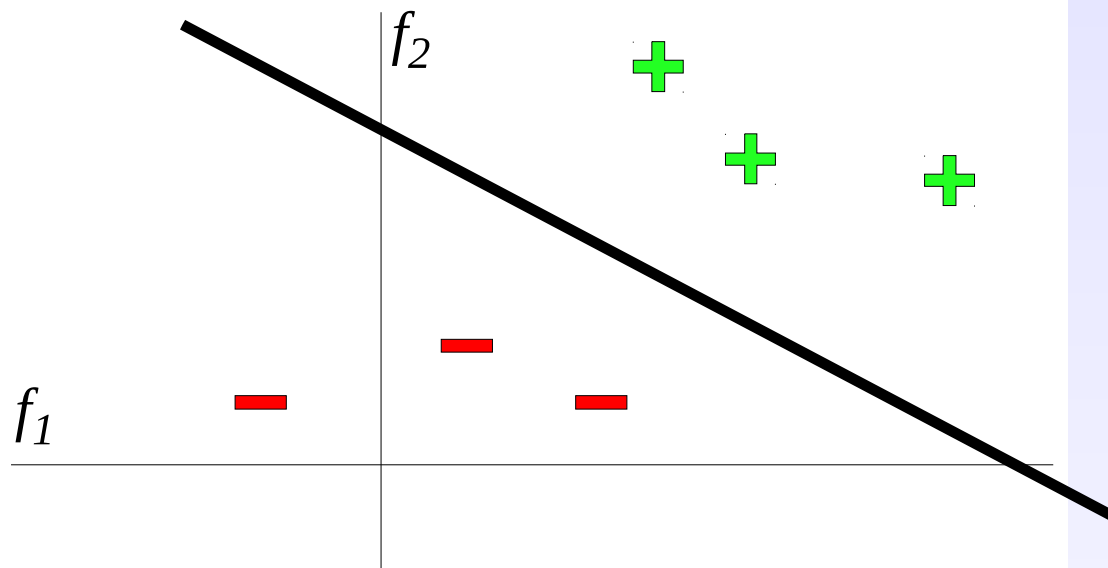
Φ

```
W=dear      :   1
W=sir       :   1
W=this      :   2
...
W=wish      :   0
...
MISSPELLED  :   2
NAMELESS    :   1
ALL_CAPS    :   0
NUM_URLS    :   0
...
```
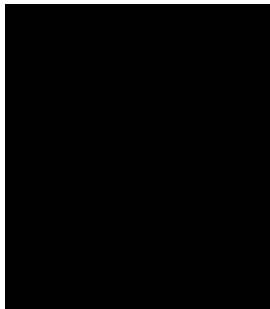
➢ Feature vectors in NLP are frequently sparse

# Linear models for binary classification

➤ Decision boundary is the set of "uncertain" points

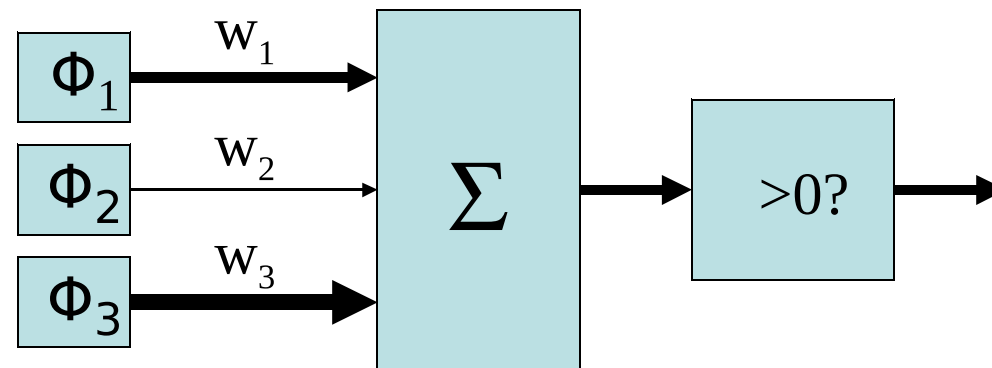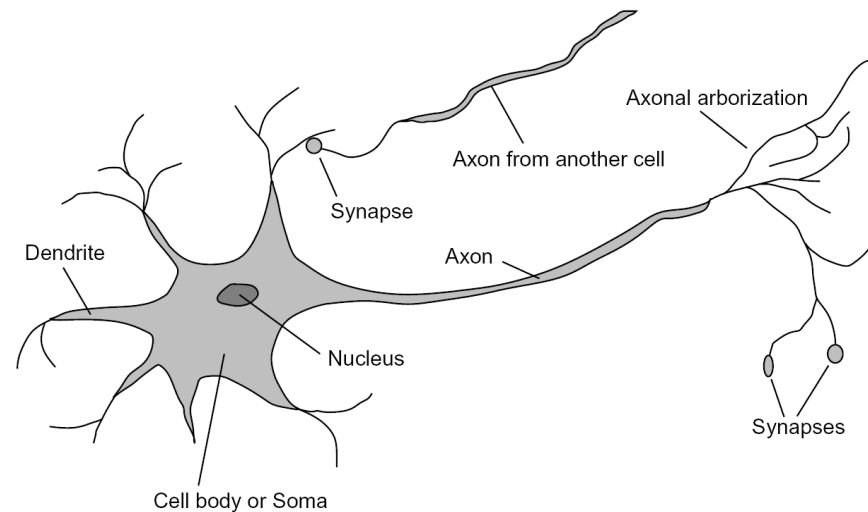➤ Linear decision boundaries are characterized by weight vectors



$$x \qquad \Phi(x) \qquad w \qquad \sum_i w_i \, \Phi_i(x)$$

"free money"

```
BIAS  :  1
free  :  1
money :  1
the   :  0
...
```

```
BIAS  : -3
free  :  4
money :  2
the   :  0
...
```

# The perceptron

- Inputs = feature values
- Params = weights
- Sum is the response

- If the response is:
  - Positive, output +1
  - Negative, output -1

___

- When training, update on errors:

$$w = w + y\,\phi(x)$$



Axonal arborization

Axon from another cell

Synapse

Dendrite

Axon

Nucleus

Synapses

Cell body or Soma

$\Phi_1$ — $w_1$ →

$\Phi_2$ — $w_2$ →

$\Phi_3$ — $w_3$ →

$\Sigma$ → >0? →

"Error" when:

$$y\,w \cdot \phi(x) \leq 0$$

# Why does that update work?

➢ When $y\,\boldsymbol{w}^{old}\cdot\phi(x)\leq 0$, update: $\boldsymbol{w}^{new}=\boldsymbol{w}^{old}+y\,\phi(x)$



$$y\,\boldsymbol{w}^{new}\,\phi(x)=y\big(\boldsymbol{w}^{old}+y\,\phi(x)\big)\phi(x)$$

$$=\underbrace{y\,\boldsymbol{w}^{old}\,\phi(x)}_{<0}+\underbrace{yy\,\phi(x)\phi(x)}_{>0}$$

# Support vector machines

➢ Explicitly optimize the ***margin***

➢ Enforce that all training points are correctly classified

$$\max_{\mathbf{w}} \quad \text{margin} \quad s.t. \quad \text{all points are correctly classified}$$

$$\max_{\mathbf{w}} \quad \text{margin} \quad s.t. \quad y_n \boldsymbol{w} \cdot \phi(x_n) \geq 1 \quad , \quad \forall n$$

$$\min_{\mathbf{w}} \quad \|\boldsymbol{w}\|^2 \quad s.t. \quad y_n \boldsymbol{w} \cdot \phi(x_n) \geq 1 \quad , \quad \forall n$$

# Support vector machines with *slack*

➢ Explicitly optimize the **margin**

➢ Allow some "noisy" points to be misclassified

$$\min_{\mathbf{w}, \xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_n \xi_n$$

$$s.t. \quad y_n \mathbf{w} \cdot \phi(x_n) + \boxed{\xi_n} \geq 1 \quad , \quad \forall n$$

$$\xi_n \geq 0 \quad , \quad \forall n$$

# Batch versus stochastic optimization

➢ Batch = read in all the data, then process it
➢ Stochastic = (roughly) process a bit at a time

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C \sum_n \xi_n$$

$$s.t. \quad y_n \mathbf{w} \cdot \phi(x_n) + \xi_n \geq 1 \\ , \ \forall n$$

$$\xi_n \geq 0 \ , \ \forall n$$

➢ For n=1..N:
  ➢ If $y_n \mathbf{w} \cdot \phi(x_n) \leq 0$
    ➢ $\mathbf{w} = \mathbf{w} + y_n \phi(x_n)$

# Stochastically optimized SVMs

SVM Objective

SOME MATH

**Implementation Note:**

Weight shrinkage is *SLOW*. Implement it lazily, at the cost of double storage.

> For n=1..N:
>> If $y_n \boldsymbol{w} \cdot \phi(x_n) \leq 1$
>>> $\boldsymbol{w} = \boldsymbol{w} + y_n \phi(x_n)$
>> $\boldsymbol{w} = \left(1 - \dfrac{1}{CN}\right) \boldsymbol{w}$

For n=1..N:
> If $y_n \boldsymbol{w} \cdot \phi(x_n) \leq 0$
>> $\boldsymbol{w} = \boldsymbol{w} + y_n \phi(x_n)$

# From Perceptron to Structured Perceptron

# Perceptron with multiple classes

- Store separate weight vector for each class $w_1, w_2, ..., w_K$

- For n=1..N:
  - Predict:
    $$\hat{y} = arg\, max_k \boldsymbol{w_k} \cdot \phi(x_n)$$

  - If $\hat{y} \neq y_n$:
    $$\boldsymbol{w_{\hat{y}}} = \boldsymbol{w_{\hat{y}}} - \phi(x_n)$$
    $$\boldsymbol{w_{y_n}} = \boldsymbol{w_{y_n}} + \phi(x_n)$$

$w_1 \bullet \Phi(x)$ biggest

$w_3 \bullet \Phi(x)$ biggest

$w_2 \bullet \Phi(x)$ biggest

**?!** Why does this do the right thing?

# **Perceptron**

|  | x | Φ(x,1) | Φ(x,2) |
|---|---|---|---|

➢ Originally:

"free money"

```
spam_BIAS  :  1
spam_free  :  1
spam_money :  1
spam_the   :  0
...
```

```
ham_BIAS  :  1
ham_free  :  1
ham_money :  1
ham_the   :  0
...
```

$w_1$

➢ For n=1..N:

➢ Predict:

$$\hat{y} = arg\,max_k \, \boldsymbol{w_k} \cdot \phi(x_n)$$

➢ If $\hat{y} \neq y_n$:

$$\boldsymbol{w_{\hat{y}}} = \boldsymbol{w_{\hat{y}}} - \phi(x_n)$$
$$\boldsymbol{w_{y_n}} = \boldsymbol{w_{y_n}} + \phi(x_n)$$

➢ For n=1..N:

➢ Predict:

$$\hat{y} = arg\,max_k \, \boldsymbol{w} \cdot \phi(x_n, k)$$

➢ If $\hat{y} \neq y_n$:

$$\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y})$$
$$+ \phi(x_n, y_n)$$

# Features for structured prediction

> Allowed to encode *anything* you want

| Pro | Md | Vb | Dt | Nn |
|-----|-----|-----|-----|-----|
| I | can | can | a | can |

$$\phi(\boldsymbol{x}, y) =$$

```
I_Pro      : 1    <s>-Pro  :  1   has_verb   :  1
can_Md     : 1    Pro-Md   :  1   has_nn_lft :  0
can_Vb     : 1    Md-Vb    :  1   has_n_lft  :  1
a_Dt       : 1    Vb-Dt    :  1   has_nn_rgt :  1
can_Nn     : 1    Dt-Nn    :  1   has_n_rgt  :  1
...               Nn-</s>  :  1   ...
                  ...
```

> Output features, Markov features, other features

# Structured perceptron

Enumeration over 1..K

Enumeration over all outputs

> For n=1..N:
>> Predict:

$$\hat{y} = arg\,max_k \, \boldsymbol{w} \cdot \phi(x_n, k)$$

>> If $\hat{y} \neq y_n$:

$$\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$$

> For n=1..N:
>> Predict:

$$\hat{y} = arg\,max_k \, \boldsymbol{w} \cdot \phi(x_n, k)$$

>> If $\hat{y} \neq y_n$:

$$\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$$

[Collins, EMNLP02]

# Argmax for sequences

> If we only have output and Markov features, we can use Viterbi algorithm:



*(plus some work to account for boundary conditions)*

# Structured perceptron as ranking

➢ For n=1..N:
  ➢ Run Viterbi: $\hat{y} = arg\,max_k\, \boldsymbol{w} \cdot \phi(x_n, k)$
  ➢ If $\hat{y} \neq y_n$: $\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$

---

➢ When does this make an update?

| Pro | Md | Vb | Dt | Nn |
|-----|-----|-----|-----|-----|
| Pro | Md | Md | Dt | Vb |
| Pro | Md | Md | Dt | Nn |
| Pro | Md | Nn | Dt | Md |
| Pro | Md | Nn | Dt | Nn |
| Pro | Md | Vb | Dt | Md |
| Pro | Md | Vb | Dt | Vb |
| I | can | can | a | can |

# From perceptron to margins

Maximize Margin

Minimize Errors

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_n \xi_n$$

$$s.t. \quad y_n \mathbf{w}\cdot\phi(x_n) + \xi_n \geq 1 \quad , \quad \forall n$$

Each point is correctly classified, modulo ξ

$$\min_{\mathbf{w},\xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_n \xi_{n,\hat{y}}$$

Response for truth

Response for other

$$s.t. \quad \mathbf{w}\cdot\phi(x_n, y_n)$$
$$-\mathbf{w}\cdot\phi(x_n, \hat{y})$$
$$+\xi_n \geq 1, \forall n, \hat{y} \neq y_n$$

Each true output is more highly ranked, modulo ξ

[Taskar+al, JMLR05; Tshochandaritis, JMLR05]

$$\min_{\mathbf{w}, \xi} \quad \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_n \xi_{n,\hat{y}}$$
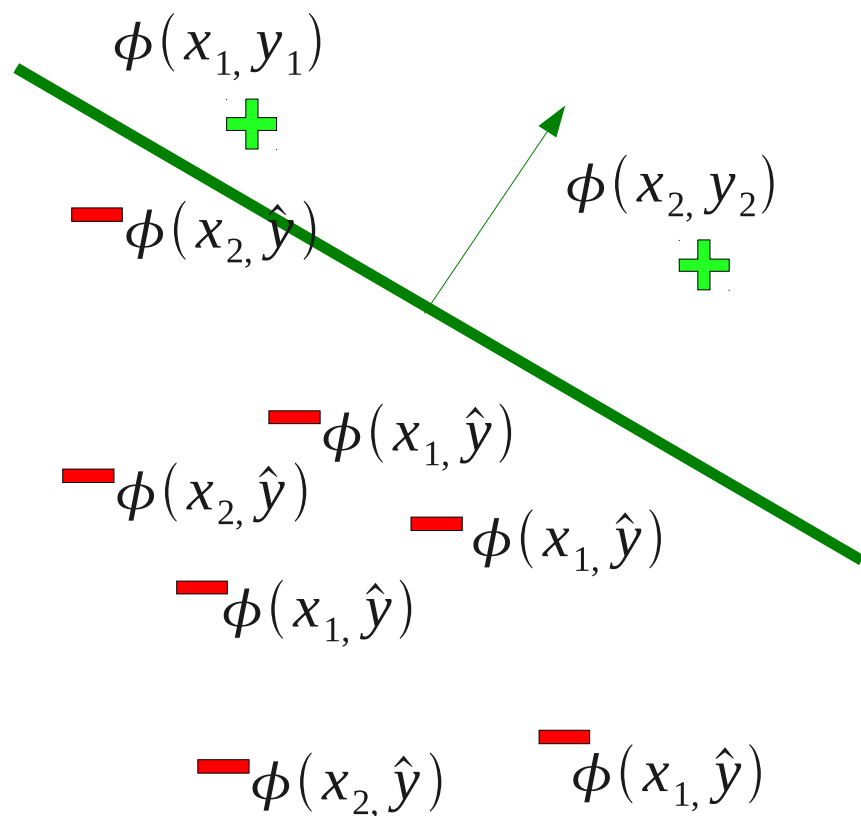
Response for truth

Response for other

$$s.t. \quad \mathbf{w}\cdot\phi(x_n, y_n)$$
$$- \mathbf{w}\cdot\phi(x_n, \hat{y})$$
$$+ \xi_n \geq 1, \forall n, \hat{y} \neq y_n$$

Each true output is more highly ranked, modulo ξ

[Taskar+al, JMLR05; Tshochandaritis, JMLR05]

➢ Some errors are worse than others…

| Pro | Md | Vb | Dt | Nn |
|-----|----|----|----|----|

*Margin of one*

| Pro | Md | Md | Dt | Vb |
|-----|----|----|----|----|
| Pro | Md | Md | Dt | Nn |
| Pro | Md | Nn | Dt | Md |
| Pro | Md | Nn | Dt | Nn |
| Pro | Md | Vb | Dt | Md |
| Pro | Md | Vb | Dt | Vb |

| I | can | can | a | can |
|---|-----|-----|---|-----|

➤ Some errors are worse than others...



*Margin of l(y,y')*

[Taskar+al, JMLR05; Tshochandaritis, JMLR05]

$\phi(x_{1,}y_1)$

$\phi(x_{2,}\hat{y})$

$\phi(x_{2,}y_2)$

$\phi(x_{1,}\hat{y})$

$\phi(x_{2,}\hat{y})$

$\phi(x_{1,}\hat{y})$

$\phi(x_{1,}\hat{y})$

$\phi(x_{2,}\hat{y})$

$\phi(x_{1,}\hat{y})$

$$\boldsymbol{w}\cdot\phi(x_n,y_n)-\boldsymbol{w}\cdot\phi(x_n,\hat{y})+\xi_n$$

$$\geq 1$$

$$\boldsymbol{w}\cdot\phi(x_n,y_n)-\boldsymbol{w}\cdot\phi(x_n,\hat{y})+\xi_n$$

$$\geq l(y_n,\hat{y})$$

[Taskar+al, JMLR05; Tshochandaritis, JMLR05]

> Add "loss" to each wrong node!

# Stochastically optimizing Markov nets

**M³N Objective**

SOME MATH

> For n=1..N:
>> Augmented Viterbi:
>> $$\hat{y} = arg\,max_k\, \boldsymbol{w} \cdot \phi(x_n, k) + l(y_n, k)$$
>> If $\hat{y} \neq y_n$:
>> $$\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$$
>> $$\boldsymbol{w} = \left(1 - \frac{1}{CN}\right)\boldsymbol{w}$$

> For n=1..N:
>> Viterbi:
>> $$\hat{y} = arg\,max_k\, \boldsymbol{w} \cdot \phi(x_n, k)$$
>> If $\hat{y} \neq y_n$:
>> $$\boldsymbol{w} = \boldsymbol{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$$

[Ratliff+al, AIStats07]

# Learning to Search

# Argmax is *hard!*

➢ Classic formulation of structured prediction:

$$score(x,y) = \quad \text{something we learn to make "good" } x,y \text{ pairs score highly}$$

➢ At test time:

$$f(x) = argmax_{y \in Y} score(x,y)$$

➢ Combinatorial optimization problem
  ➢ Efficient only in very limiting cases
  ➢ Solved by heuristic search: beam + A* + local search

S

Train a classifier to make decisions

**Right**

**Left**

VP

**Right**

NP

NP

**Unary**

**Left**

**Up**

**Left**

**Right**

I / Pro

can / Md

can / Vb

a / Dt

can / Nn

[Magerman, ACL95]

# Incremental parsing, mid 2000s style

# Learning to beam-search

**Beam-search**

- For n=1..N:
  - ~~Viterbi:~~
    $$\hat{y} = \arg\max_k \mathbf{w} \cdot \phi(x_n, k)$$
  - If $\hat{y} \neq y_n$

$$\mathbf{w} = \mathbf{w} - \phi(x_n, \hat{y}) + \phi(x_n, y_n)$$

S

VP

NP

NP

I / Pro    can / Md    can / Vb    a / Dt    can / Nn

# Learning to beam-search

S

VP

> For n=1..N:
>> Run beam search until truth falls out of beam
>> Update weights immediately!

NP

NP

I / Pro

can / Md

can / Vb

a / Dt

can / Nn

# Learning to beam-search

S

VP

For n=1..N:
- Run beam search until truth falls out of beam
- Update weights immediately!
- Restart at truth

NP

NP

I / Pro      can / Md      can / Vb      a / Dt      can / Nn

# Incremental parsing results



Figure legend:
- No early update, no repeated use of examples
- Early update, no repeated use of examples
- Early update, repeated use of examples

Y-axis: F-measure parsing accuracy (82 to 88)

X-axis: Number of passes over training data (1 to 6)

[Collins+Roark, ACL04]

# Generic Search Formulation

- Search Problem:
  - Search space
  - Operators
  - Goal-test function
  - Path-cost function

- Search Variable:
  - Enqueue function

Varying the **Enqueue** function can give us DFS, BFS, beam search, A* search, etc...

- nodes := MakeQueue(S0)

- **while** nodes is not empty
  - node := RemoveFront(nodes)
  - **if** node is a goal state **return** node
  - next := Operators(node)
  - nodes := Enqueue(nodes, next)

- **fail**

[D+Marcu, ICML05; Xu+al, JMLR09]

> nodes := MakeQueue(S0)
>
> **while** nodes is not empty
>> node := RemoveFront(nodes)
>>
>> **if** none of {node} ∪ nodes is y-good **or** node is a goal **&** not y-good

*Monotonicity*: for any node, we can tell if it can lead to the correct solution or not

If we erred...

Where should we have gone?

>>> sibs := siblings(node, y)
>>>
>>> w := update(w, x, sibs, {node} ∪ nodes)
>>>
>>> nodes := MakeQueue(sibs)

Update our weights based on the good and the bad choices

>> **else**

Continue search...

>>> **if** node is a goal state **return** w
>>>
>>> next := Operators(node)
>>>
>>> nodes := Enqueue(nodes, next)

[D+Marcu, ICML05; Xu+al, JMLR09]

# Search-based Margin

- The *margin* is the amount by which we are correct:



$$\bar{u}^T \Phi(x, g_1)$$

$$\bar{u}^T \Phi(x, g_2)$$

$$\bar{u}^T \Phi(x, b_1)$$

$$\bar{u}^T \Phi(x, b_2)$$

$\bar{u}$

$\gamma$

- Note that the *margin* and hence *linear separability* is also a function of the *search algorithm!*

# Syntactic chunking Results



Hal Daumé III (me@hal3.name)  SP2IRL @ ACL2010

# Tagging+chunking results



Hal Daumé III (me@hal3.name) SP2IRL @ ACL2010

# Variations on a beam

- Observation:
  - We needn't use the same beam size for training and decoding
  - Varying these values independently yields:

<div align="center">

**Decoding Beam**

| Training Beam | 1 | 5 | 10 | 25 | 50 |
|---|---|---|---|---|---|
| 1 | 93.9 | 92.8 | 91.9 | 91.3 | 90.9 |
| 5 | 90.5 | 94.3 | 94.4 | 94.1 | 94.1 |
| 10 | 89.5 | 94.3 | 94.4 | 94.2 | 94.2 |
| 25 | 88.7 | 94.2 | 94.5 | 94.3 | 94.3 |
| 50 | 88.4 | 94.2 | 94.4 | 94.2 | 94.4 |

</div>

[D+Marcu, ICML05; Xu+al, JMLR09]

# What if our model sucks?

➢ Sometimes our model *cannot* produce the "correct" output

  ➢ canonical example: machine translation



"Bold" update

Current Hypothesis

Reference

**Good Outputs**

N-best list *or* "optimal decoding" *or* ...

**Model Outputs**

"Local" update

Best achievable output

[Och, ACL03; Liang+al, ACL06]

# Local versus bold updating...



Machine Translation Performance (Bleu)

Legend:
- Bold
- Local
- Pharoah

X-axis: Monotonic, Distortion
Y-axis: 32.5, 33, 33.5, 34, 34.5, 35, 35.5

# Integrating search and learning

**Input:** Le homme mange l' croissant.
**Output:** The man ate a croissant.

Hyp: The man ate
Cov: Le homme mange l' croissant.

Classifier 'h'

Hyp: The man ate a croissant
Cov: Le homme mange l' croissant.

Hyp: The man ate a fox
Cov: Le homme mange l' croissant.

Hyp: The man ate happy
Cov: Le homme mange l' croissant.

Hyp: The man ate a
Cov: Le homme mange l' croissant.

Hyp: The man ate a
Cov: Le homme mange l' croissant.

[D+Marcu, ICML05; D+Langford+Marcu, MLJ09]

# Reducing search to classification

- Natural chicken and egg problem:
  - Want *h* to get low expected future loss
  - … on future decisions made by *h*
  - … and starting from states visited by *h*
- Iterative solution

$h^{(t-1)}$

**Input:** Le homme mange l' croissant.
**Output:** The man ate a croissant.

Hyp: The man ate
Cov: Le homme mange
l' croissant.

$h^{(t)}$

Hyp: The man ate a croissant
Cov: Le homme mange
l' croissant.
→ Loss = 0

Hyp: The man ate a fox
Cov: Le homme mange
l' croissant.
→ Loss = 1.8

Hyp: The man ate happy
Cov: Le homme mange
l' croissant.
→ Loss = 1.2

Hyp: The man ate a
Cov: Le homme mange
l' croissant.
→ Loss = 0.5

Hyp: The man ate a
Cov: Le homme mange
l' croissant.
→ Loss = 0

[D+Langford+Marcu, MLJ09]

# Theoretical results

**Theorem:** After $2T^3 \ln T$ iterations, the loss of the learned policy is bounded as follows:

$$L(h) \le L(h_0) + 2\,T \ln T \, l_{avg} + (1 + \ln T)\frac{c_{max}}{T}$$

Loss of the optimal policy

Average multiclass classification loss

Worst case per-step loss

# Example task: summarization

*That's perfect!*

Standard approach is sentence extraction, but that is often deemed to "coarse" to produce good, very short summaries. We wish to also drop words and phrases => document compression

Argentina was still obsessed with the Falkland Islands even in 1994, 12 years after its defeat in the 74-day war with Britain. The country's overriding foreign policy aim continued to be winning sovereignty over the islands.

The Falkland islands war, in 1982, was fought between Britain and Argentina.

[D+Langford+Marcu, MLJ09]

# Structure of search

Argentina was still obsessed with the Falkland Islands even in 1994, 12 years after its defeat in the 74-day war with Britain. The country's overriding foreign policy aim continued to be winning sovereignty over the islands.

➢ Lay sentences out sequentially
➢ Generate a dependency parse of each sentence
➢ Mark each root as a frontier node
➢ Repeat:
  ➢ Choose a frontier node node to add to the summary
  ➢ Add all its children to the frontier
  ➢ Finish when we have enough words



S1        S2        S3        S4        S5    ...    Sn

⬤ = frontier node     ⬤ = summary node

# Example output (40 word limit)

**Sentence Extraction + Compression:**

**+13** Argentina and Britain announced an agreement, nearly eight years after they fought a 74-day war a populated archipelago off Argentina's coast. Argentina gets out the red carpet, official royal visitor since the end of the Falklands war in 1982.

**Vine Growth (Searn):**

**+24** Argentina and Britain announced to restore full ties, eight years after they fought a 74-day war over the Falkland islands. Britain invited Argentina's minister Cavallo to London in 1992 in the first official visit since the Falklands war in 1982.

| | |
|---|---|
| 6 Diplomatic ties restored | 3 Falkland war was in 1982 |
| 5 Major cabinet member visits | 3 Cavallo visited UK |
| 5 Exchanges were in 1992 | 2 War was 74-days long |
| 3 War between Britain and Argentina | |

# Perceptron vs. LaSO vs. Searn

● Incremental perceptron

● LaSO

● Searn

✗ Un-learnable decision

# Take-home messages

➤ If you can predict (ie., solve argmax) you can learn (use structured perceptron)

➤ If you can do loss-augmented search, you can do max margin (add two lines of code to perceptron)

➤ If you can do beam search, you can learn using LaSO (with no loss function)

➤ If you can do beam search, you can learn using Searn (with any loss function)

# Coffee Break!!!

Hal Daumé III (me@hal3.name)                          SP2IRL @ ACL2010

# Refresher on Reinforcement Learning

# Reinforcement learning

- Basic idea:
  - Receive feedback in the form of <span style="color:red">rewards</span>
  - Agent's utility is defined by the reward function
  - Must learn to act to <span style="color:red">maximize expected rewards</span>
  - <span style="color:red">Change the rewards, change the learned behavior</span>

- Examples:
  - Playing a game, reward at the end for outcome
  - Vacuuming, reward for each piece of dirt picked up
  - Driving a taxi, reward for each passenger delivered

# Markov decision processes

What are the values (expected future rewards) of states and actions?

$V(s)* = $ 30

$Q(s,a1)* = $ 30

0.6

0.4

5

1

$Q(s,a2)* = $ 23

1.0

3

$Q(s,a3)* = $ 17

0.1

0.9

2

4

# Markov Decision Processes

> An MDP is defined by:
>> A set of states s ∈ S
>> A set of actions a ∈ A
>> A transition function T(s,a,s')
>>> Prob that a from s leads to s
>>> i.e., P(s' | s,a)
>>> Also called the model
>> A reward function R(s, a, s')
>>> Sometimes just R(s) or R(s')
>> A start state (or distribution)
>> Maybe a terminal state

> MDPs are a family of non-deterministic search problems
> Total utility is one of:

$$\sum_t r_t \quad \text{or} \quad \sum_t \gamma^t r_t$$

# Solving MDPs

➢ In deterministic single-agent search problem, want an optimal plan, or sequence of actions, from start to a goal

➢ In an MDP, we want an optimal policy π(s)

  ➢ A policy gives an action for each state
  ➢ Optimal policy maximizes expected if followed
  ➢ Defines a reflex agent

Optimal policy when
$R(s, a, s') = -0.04$ for
all non-terminals s

# Example Optimal Policies



R(s) = -0.01

R(s) = -0.03

R(s) = -0.4

R(s) = -2.0

Hal Daumé III (me@hal3.name) SP2IRL @ ACL2010

# Optimal Utilities

- Fundamental operation: compute the optimal utilities of states s (all at once)

- Why?  Optimal values define optimal policies!

- Define the utility of a state s:
  $V^*(s)$ = expected return starting in s and acting optimally

- Define the utility of a q-state (s,a):
  $Q^*(s,a)$ = expected return starting in s, taking action a and thereafter acting optimally

- Define the optimal policy:
  $\pi^*(s)$ = optimal action from state s

# The Bellman Equations

➢ Definition of utility leads to a simple one-step lookahead relationship amongst optimal utility values:

Optimal rewards = maximize over first action and then follow optimal policy

➢ Formally:

$$V^*(s) = \max_a Q^*(s, a)$$

$$V^*(s) = \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V^*(s') \right]$$

# Solving MDPs / memoized recursion

➢ Recurrences:

$$V_0^*(s) = 0$$

$$V_i^*(s) = \max_a Q_i^*(s, a)$$

$$Q_i^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_{i-1}^*(s') \right]$$

$$\pi_i(s) = \arg \max_a Q_i^*(s, a)$$

➢ Cache all function call results so you never repeat work

➢ What happened to the evaluation function?

# Q-Value Iteration

- Value iteration: iterate approx optimal values
  - Start with $V_0^*(s) = 0$, which we know is right (why?)
  - Given $V_i^*$, calculate the values for all states for depth i+1:

$$V_{i+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma V_i(s') \right]$$

- But Q-values are more useful!
  - Start with $Q_0^*(s,a) = 0$, which we know is right (why?)
  - Given $Q_i^*$, calculate the q-values for all q-states for depth i+1:

$$Q_{i+1}(s, a) \leftarrow \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q_i(s', a') \right]$$

# RL = Unknown MDPs

➢ If we *knew* the MDP (i.e., the reward function and transition function):

  ➢ Value iteration leads to optimal values

  ➢ Q-value iteration leads to optimal Q-values

  ➢ Will always converge to the truth

➢ Reinforcement learning is what we do when we *do not know* the MDP

  ➢ All we observe is a *trajectory*

  ➢ $(s_1, a_1, r_1, \quad s_2, a_2, r_2, \quad s_3, a_3, r_3, \quad ...)$

# Q-Learning

- Learn Q*(s,a) values
  - Receive a sample **(s,a,s',r)**
  - Consider your old estimate: $Q(s, a)$
  - Consider your new sample estimate:

$$Q^*(s, a) = \sum_{s'} T(s, a, s') \left[ R(s, a, s') + \gamma \max_{a'} Q^*(s', a') \right]$$

  - Incorporate the new estimate into a running average:

$$sample = R(s, a, s') + \gamma \max_{a'} Q(s', a')$$

$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + (\alpha) \left[ sample \right]$$

# Exploration / Exploitation

- Several schemes for forcing exploration
  - Simplest: random actions ($\varepsilon$ greedy)
    - Every time step, flip a coin
    - With probability $\varepsilon$, act randomly
    - With probability 1-$\varepsilon$, act according to current policy

  - Problems with random actions?
    - You do explore the space, but keep thrashing around once learning is done
    - One solution: lower $\varepsilon$ over time
    - Another solution: exploration functions

Hal Daumé III (me@hal3.name) SP2IRL @ ACL2010

# Q-Learning

> In realistic situations, we cannot possibly learn about every single state!
>> Too many states to visit them all in training
>> Too many states to hold the q-tables in memory

> Instead, we want to generalize:
>> Learn about some small number of training states from experience
>> Generalize that experience to new, similar states:

███████████████████████████████████

> Very simple stochastic updates:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \, [error]$$

$$w_i \leftarrow w_i + \alpha \, [error] \, f_i(s,a)$$

# Inverse RL *and* Apprenticeship Learning

# Inverse RL: Task

> Given:

>> measurements of an agent's behavior over time, in a variety of circumstances

>> if needed, measurements of the sensory inputs to that agent

>> if available, a model of the environment.

> Determine: the reward function being optimized

> Proposed by [Kalman68]

> First solution, by [Boyd94]

# Why inverse RL?

- Computational models for animal learning
  - "In examining animal and human behavior we must consider the reward function as an unknown to be ascertained through empirical investigation."

- Agent construction
  - "An agent designer [...] may only have a very rough idea of the reward function whose optimization would generate 'desirable' behavior."
  - eg., "Driving well"

- Multi-agent systems and mechanism design
  - learning opponents' reward functions that guide their actions to devise strategies against them

# IRL from Sample Trajectories

➢ Optimal policy available through trajectories (eg., driving a car)

➢ Want to find *Reward* function that makes this policy look *as good as possible*

➢ Write $R_w(s) = w\,\phi(s)$ so the reward is linear

and $V_w^\pi(s_0)$ be the value of the starting state

$$\max_{\mathbf{w}} \sum_{k=1}^{K} f\left(V_w^{\pi^*}(s_0) - V_w^{\pi_k}(s_0)\right)$$

How good does the "optimal policy" look?

How good does the some other policy look?

[Ng+Russell, ICML00]

# Apprenticeship Learning via IRL

➢ For  t = 1,2,…

  ➢ Inverse RL step:
    Estimate expert's reward function $R(s) = w^\top \phi(s)$ such that under R(s) the expert performs better than all previously found policies $\{\pi_i\}$.

  ➢ RL step:
    Compute optimal policy $\pi_t$ for the estimated reward w

[Abbeel+Ng, ICML04]

# Car Driving Experiment

- No explicit reward function at all!

- Expert demonstrates proper policy via 2 min. of driving time on simulator (1200 data points).

- 5 different "driver types" tried.

- Features: which lane the car is in, distance to closest car in current lane.

- Algorithm run for 30 iterations, policy hand-picked.

- Movie Time!  (Expert left, IRL right)

[Abbeel+Ng, ICML04]

# "Nice" driver

Hal Daumé III (me@hal3.name)

# "Evil" driver



Hal Daumé III (me@hal3.name)    SP2IRL @ ACL2010

# Maxent IRL

Distribution over trajectories:
$$P(\zeta)$$

Match the reward of observed behavior:

$$\sum_\zeta P(\zeta)\, f_\zeta = f_{\textbf{dem}}$$

Maximizing the **causal entropy** over trajectories given stochastic outcomes:

$$\max \textbf{H}(\textbf{P}(\zeta)||\textbf{O})$$

(Condition on random uncontrolled outcomes, but only **after** they happen)

**As uniform as possible**

# Data collection



**Length Speed Road Type Lanes**

**Accidents Construction Congestion Time of day**

**25 Taxi Drivers**

**Over 100,000 miles**

[Ziebart+al, AAAI08]

# Predicting destinations....

# Inverse Optimal Control

# Planning as structured prediction

Hal Daumé III (me@hal3.name)

SP2IRL @ ACL2010

[Ratliff+al, NIPS05]

# Maximum margin planning

> Let µ(s,a) denote the probability of reaching q-state (s,a) under current model w

$$\max_{\mathbf{w}} \text{ margin } \quad s.t. \quad \text{planner run with w yields human output}$$

Q-state visitation frequency by human

$$\min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w}\|^2 \quad s.t. \quad \begin{array}{l} \mu(s,a)\mathbf{w}\cdot\phi(x_n,s,a) \\ -\hat{\mu}(s,a)\mathbf{w}\cdot\phi(x_n,s,a) \geq 1 \\ , \quad \forall n,s,a \end{array}$$

Q-state visitation frequency by planner

All trajectories, and all q-states

# Optimizing MMP

M³N Objective

SOME MATH



- For n=1..N:
  - Augmented planning:
    Run A* on current (augmented) cost map
    to get q-state visitation frequencies $\mu(s,a)$

  - Update: $w = w + \sum_s \sum_a \left[\hat{\mu}(s,a) - \mu(s,a)\right] \phi(x_n, s, a)$

  - Shrink: $w = \left(1 - \dfrac{1}{CN}\right) w$

# Maximum margin planning movies

[Ratliff+al, NIPS05]

# Parsing via inverse optimal control

- State space = all partial parse trees over the full sentence labeled "S"
- Actions: take a partial parse and split it anywhere in the middle
- Transitions: obvious
- Terminal states: when there are no actions left
- Reward: parse score at completion

[Neu+Szepevari, MLJ09]

# Parsing via inverse optimal control



Small      Medium      Large

- ■ Maximum Likelihood
- ■ Projection
- ■ Perceptron
- ■ Apprenticeship Learning
- ■ Maximum Margin
- ■ Maximum Entropy
- ■ Policy Matching

[Neu+Szepevari, MLJ09]

# Learning to Search

$(f^p, +1)$

$f_2$

$(f^d, -1)$

$f_1$

[Ratliff+al, AutRobots09]

# Learch

Until converged do:

Initialize modification set to empty

For each example, add cost function modifications:

Make loss-augmented prediction using current cost

Update data set: Label desired feature vector as -1 and predicted feature vector as +1

Generalize using a least-squares regression

Add it to the current cost function

Hal Daumé III (me@hal3.name)

# **Discussion**

Hal Daumé III (me@hal3.name)                    SP2IRL @ ACL2010

# Relationship between SP and IRL

➢ Formally, they're (nearly) the same problem

   ➢ See humans performing some task

   ➢ Define some loss function

   ➢ Try to mimic the humans

➢ Difference is in philosophy:

   ➢ (I)RL has little notion of beam search or dynamic programming

   ➢ SP doesn't think about separating reward estimation from solving the prediction problem

   ➢ (I)RL has to deal with stochastiticity in MDPs

# Important Concepts

➢ Search and loss-augmented search for margin-based methods

➢ Bold versus local updates for approximate search

➢ Training on-path versus off-path

➢ Stochastic versus deterministic worlds

➢ Q-states / values

➢ Learning reward functions vs. matching behavior

# Hal's Wager

- ➤ Give me a structured prediction problem where:
    - ➤ Annotations are at the lexical level
    - ➤ Humans can do the annotation with reasonable agreement
    - ➤ You give me a few thousand labeled sentences

- ➤ Then I can learn reasonably well...
    - ➤ ...using one of the algorithms we talked about

- ➤ Why do I say this?
    - ➤ Lots of positive experience
    - ➤ I'm an optimist
    - ➤ I want your *counter-examples!*

# Open problems

- How to do SP when argmax is intractable….
  - Bad: simple algorithms diverge  [Kulesza+Pereira, NIPS07]
  - Good: some work well  [Finley+Joachims, ICML08]
  - And you can make it fast!  [Meshi+al, ICML10]

- How to do SP with delayed feedback (credit assignment)
  - Kinda just works sometimes  [D, ICML09; Chang+al, ICML10]
  - Generic RL also works  [Branavan+al, ACL09; Liang+al, ACL09]

- What role does structure actually play?
  - Little: only constraints outputs  [Punyakanok+al, IJCAI05]
  - Little: only introduces non-linearities  [Liang+al, ICML08]
  - Lots: ???

# Things I have no idea how to solve...

**all :** (a → Bool) → [a] → Bool

Applied to a predicate and a list, returns `True` if all elements
of the list satisfy the predicate, and `False` otherwise.

**all** p
**all** p
  **if** p
     **the**
     **els**

```
%module main:MyPrelude
   %data main:MyPrelude.MyList aadj =
      {main:MyPrelude.Nil;
       main:MyPrelude.Cons aadj ((main:MyPrelude.MyList aadj))};
   %rec
   {main:MyPrelude.myzuall :: %forall tadA . (tadA ->
                                                ghczmprim:GHCziBool.Bool)
                                             ->
                                             (main:MyPrelude.MyList tadA) ->
                                             ghczmprim:GHCziBool.Bool =
      \ @ tadA
        (padk::tadA -> ghczmprim:GHCziBool.Bool)
        (dsddE::(main:MyPrelude.MyList tadA)) ->
          %case ghczmprim:GHCziBool.Bool dsddE
          %of (wildB1::(main:MyPrelude.MyList tadA))
             {main:MyPrelude.Nil ->
                 ghczmprim:GHCziBool.True;
              main:MyPrelude.Cons
               (xadm::tadA) (xsadn::(main:MyPrelude.MyList tadA)) ->
                 %case ghczmprim:GHCziBool.Bool (padk xadm)
                 %of (wild1Xc::ghczmprim:GHCziBool.Bool)
                    {ghczmprim:GHCziBool.False ->
                        ghczmprim:GHCziBool.False;
                     ghczmprim:GHCziBool.True ->
                        main:MyPrelude.myzuall @ tadA padk xsadn}}};
```

(s1) A father had a family of sons who were perpetually quarreling among themselves. (s2) When he failed to heal their disputes by his exhortations, he determined to give them a practical illustration of the evils of disunion; and for this purpose he one day told them to bring him a bundle of sticks. (s3) When they had done so, he placed the faggot into the hands of [...] them to break it i[...] strength, and we[...] the faggot, took t[...] again put them i[...] them easily. (s6) [...] "My sons, if you [...] other, you will be[...] of your enemies; [...] you will be broke[...]



| Father | | Sons |
|---|---|---|
| | shared......... $-_{s1}$(quarreling)$_{a1}$ | |
| $(annoyed)_{a2}$ .......... | | |
| $M_{s2}$ $(stop\ quarreling)_{a3}$ | | |
| $M_{s2}$ $(exhortations)_{a4}$ | | |
| $-_{s2}$ $(exhortations\ fail)_{a5}$ | | |
| $M_{s2}$ $(teach\ lesson)_{a6}$ | | |
| $M_{s2}$ $(get\ sticks\ \&\ break)_{a7}$ ....request | $M_{s2}(get\ sticks\ \&\ break)_{a8}$ | |
| | shared ..... $-_{s4}$(cannot break sticks)$_{a9}$ | |
| $+_{s4}$ $(cannot\ break\ sticks)_{a10}$ | | |
| $M_{s5}$ $(bundle\ \&\ break)_{a11}$....request | $M_{s5}(bundle\ \&\ break)_{a12}$ | |
| | shared ......... $+_{s5}$(break sticks)$_{a13}$ | |
| $+_{s5}$ $(break\ sticks)_{a14}$ | | |
| $+_{s5}$ $(lesson\ succeeds)_{a15}$ | | |

# Software

- Sequence labeling
    - Mallet      http://mallet.cs.umass.edu
    - CRF++       http://crfpp.sourceforge.net

- Search-based structured prediction
    - LaSO        http://hal3.name/TagChunk
    - Searn       http://hal3.name/searn

- Higher-level "feature template" approaches
    - Alchemy     http://alchemy.cs.washington.edu
    - Factorie    http://code.google.com/p/factorie

# Summary

> Structured prediction is *easy* if you can do argmax search (esp. loss-augmented!)

> Label-bias can kill you, so iterate (Searn)

> Stochastic worlds modeled by MDPs

> IRL is all about learning reward functions

> IRL has fewer assumptions

>> More general

>> Less likely to work on easy problems

> We're a long way from a complete solution

> Hal's wager: we can learn pretty much anything

# Thanks! Questions?

# References

See also:

http://www.cs.utah.edu/~suresh/mediawiki/index.php/MLRG
http://braque.cc/ShowChannel?handle=P5BVAC34

Hal Daumé III (me@hal3.name) SP2IRL @ ACL2010

# Stuff we talked about explicitly

➤ *Apprenticeship learning via inverse reinforcement learning, P. Abbeel and A. Ng. ICML, 2004.*

➤ *Incremental parsing with the Perceptron algorithm.* M. Collins and B. Roark. ACL 2004.

➤ *Discriminative Training Methods for Hidden Markov Models: Theory and Experiments with Perceptron Algorithms.* M. Collins. EMNLP 2002.

➤ *Search-based Structured Prediction.* H. Daumé III, J. Langford and D. Marcu. Machine Learning, 2009.

➤ *Learning as Search Optimization: Approximate Large Margin Methods for Structured Prediction.* H. Daumé III and D. Marcu. ICML, *2005.*

➤ *An End-to-end Discriminative Approach to Machine Translation.* P. Liang, A. Bouchard-Côté, D. Klein, B. Taskar. ACL 2006.

➤ *Statistical Decision-Tree Models for Parsing. D. Magerman. ACL 1995.*

➤ *Training Parsers by Inverse Reinforcement Learning. G. Neu and Cs. Szepesvári. Machine Learning 77, 2009.*

➤ *Algorithms for inverse reinforcement learning, A. Ng and A. Russell. ICML, 2000.*

➤ *(Online) Subgradient Methods for Structured Prediction.* N. Ratliff, J. Bagnell, and M. Zinkevich. AIStats 2007.

➤ *Maximum margin planning.* N. Ratliff, J. Bagnell and M. Zinkevich. ICML, 2006.

➤ *Learning to search: Functional gradient techniques for imitation learning.* N. Ratliff, D. Silver, and J. Bagnell. Autonomous Robots, Vol. 27, No. 1, July, 2009.

➤ *Max-Margin Markov Networks.* B. Taskar, C. Guestrin, V. Chatalbashev and D. Koller. JMLR 2005.

➤ *Large Margin Methods for Structured and Interdependent Output Variables.* I. Tsochantaridis, T. Joachims, T. Hofmann, and Y. Altun. JMLR 2005.

➤ *Learning Linear Ranking Functions for Beam Search with Application to Planning.* Y. Xu, A. Fern, and S. Yoon. JMLR 2009.

➤ *Maximum Entropy Inverse Reinforcement Learning.* B. Ziebart, A. Maas, J. Bagnell, and A. Dey. AAAI 2008.

# Other good stuff

- *Reinforcement learning for mapping instructions to actions.* S.R.K. Branavan, H. Chen, L. Zettlemoyer and R. Barzilay. ACL, 2009.
- *Driving semantic parsing from the world's response.* J. Clarke, D. Goldwasser, M.-W. Chang, D. Roth. CoNLL 2010.
- *New Ranking Algorithms for Parsing and Tagging: Kernels over Discrete Structures, and the Voted Perceptron.* M.Collins and N. Duffy. ACL 2002.
- *Unsupervised Search-based Structured Prediction.* H. Daumé III. ICML 2009.
- *Training structural SVMs when exact inference is intractable.* T. Finley and T. Joachims. ICML, 2008.
- *Structured learning with approximate inference.* A. Kulesza and F. Pereira. NIPS, 2007.
- *Conditional random fields: Probabilistic models for segmenting and labeling sequence data.* J. Lafferty, A. McCallum, F. Pereira. ICML 2001.
- *Structure compilation: trading structure for features.* P. Liang, H. Daume, D. Klein. ICML 2008.
- *Learning semantic correspondences with less supervision.* P. Liang, M. Jordan and D. Klein. ACL, 2009.
- *Generalization Bounds and Consistency for Structured Labeling.* D. McAllester. In *Predicting Structured Data, 2007.*
- *Maximum entropy Markov models for information extraction and segmentation.* A. McCallum, D. Freitag, F. Pereira. ICML 2000.
- *FACTORIE: Efficient Probabilistic Programming for Relational Factor Graphs via Imperative Declarations of Structure, Inference and Learning.* A. McCallum, K. Rohanemanesh, M. Wick, K. Schultz, S. Singh. NIPS Workshop on Probabilistic Programming, 2008
- *Learning efficiently with approximate inference via dual losses.* O. Meshi, D. Sontag, T. Jaakkola, A. Globerson. ICML 2010.
- *Learning and inference over constrained output.* V. Punyakanok, D. Roth, W. Yih, D. Zimak. IJCAI, 2005.
- *Boosting Structured Prediction for Imitation Learning.* N. Ratliff, D. Bradley, J. Bagnell, and J. Chestnutt. NIPS 2007.
- *Efficient Reductions for Imitation Learning.* S. Ross and J. Bagnell. AISTATS, 2010.
- *Kernel Dependency Estimation.* J. Weston, O. Chapelle, A. Elisseeff, B. Schoelkopf and V. Vapnik. NIPS 2002.