# Algorithms for NLP



## Classification III

Taylor Berg-Kirkpatrick – CMU

Slides: Dan Klein – UC Berkeley

# The Perceptron, Again

- Start with zero weights
- Visit training instances one by one
    - Try to classify

$$\widehat{\mathbf{y}} = \arg\max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})$$

    - If correct, no change!
    - If wrong: adjust weights

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{f}_i(\mathbf{y}_i^*)$$
$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{f}_i(\widehat{\mathbf{y}})$$

$$\mathbf{w} \leftarrow \mathbf{w} + (\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\widehat{\mathbf{y}}))$$

$$\mathbf{w} \leftarrow \mathbf{w} + \boxed{\Delta_i(\widehat{\mathbf{y}})} \quad \textit{mistake vectors}$$

# Perceptron Weights

$$\mathbf{w} \leftarrow \mathbf{w} + \Delta_i(\mathbf{y})$$

- **What is the final value of w?**
  - Can it be an arbitrary real vector?
  - No! It's built by adding up feature vectors (mistake vectors).

$$\mathbf{w} = \Delta_i(\mathbf{y}) + \Delta_{i'}(\mathbf{y'}) + \cdots$$

$$\mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \Delta_i(\mathbf{y}) \qquad \textit{mistake counts}$$

- **Can reconstruct weight vectors (the primal representation) from update counts (the dual representation) for each i**

$$\alpha_i = \langle \alpha_i(\mathbf{y_1}) \ \ \alpha_i(\mathbf{y_2}) \ \ \ldots \ \ \alpha_i(\mathbf{y_n}) \rangle$$

# Dual Perceptron

- Track mistake counts rather than weights

$$\mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \Delta_i(\mathbf{y})$$

- Start with zero counts ($\alpha$)
- For each instance x
  - Try to classify

$$\hat{\mathbf{y}} = \arg\max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mathbf{w}^\top \mathbf{f}(\mathbf{y})$$

$$\hat{\mathbf{y}} = \arg\max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x}_i)} \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \Delta_{i'}(\mathbf{y}')^\top \mathbf{f}_i(\mathbf{y})$$

  - If correct, no change!
  - If wrong: raise the mistake count for this example and prediction

$$\alpha_i(\hat{\mathbf{y}}) \leftarrow \alpha_i(\hat{\mathbf{y}}) + 1 \qquad \mathbf{w} \leftarrow \mathbf{w} + \Delta_i(\hat{\mathbf{y}})$$

# Dual / Kernelized Perceptron

- How to classify an example x?

$$score(\mathbf{y}) = \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) = \left( \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \Delta_{i'}(\mathbf{y}') \right)^\top \mathbf{f}_i(\mathbf{y})$$

$$= \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \left( \Delta_{i'}(\mathbf{y}')^\top \mathbf{f}_i(\mathbf{y}) \right)$$

$$= \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \left( \mathbf{f}_{i'}(\mathbf{y}_{i'}^*)^\top \mathbf{f}_i(\mathbf{y}) - \mathbf{f}_{i'}(\mathbf{y}')^\top \mathbf{f}_i(\mathbf{y}) \right)$$

$$= \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \left( K(\mathbf{y}_{i'}^*, \mathbf{y}) - K(\mathbf{y}', \mathbf{y}) \right)$$

- If someone tells us the value of K for each pair of candidates, never need to build the weight vectors

# Issues with Dual Perceptron

- Problem: to score each candidate, we may have to compare to *all* training candidates

$$score(\mathbf{y}) = \sum_{i',\mathbf{y}'} \alpha_{i'}(\mathbf{y}') \left( K(\mathbf{y}_{i'}^*, \mathbf{y}) - K(\mathbf{y}', \mathbf{y}) \right)$$

  - Very, very slow compared to primal dot product!
  - One bright spot: for perceptron, only need to consider candidates we made mistakes on during training
  - Slightly better for SVMs where the alphas are (in theory) sparse

- This problem is serious: fully dual methods (including kernel methods) tend to be extraordinarily slow
- Of course, we can (so far) also accumulate our weights as we go...

# Kernels: Who Cares?

- So far: a very strange way of doing a very simple calculation

- "Kernel trick": we can substitute any* similarity function in place of the dot product

- Lets us learn new kinds of hypotheses

* Fine print: if your kernel doesn't satisfy certain technical requirements, lots of proofs break. E.g. convergence, mistake bounds. In practice, illegal kernels *sometimes* work (but not always).

# Some Kernels

- Kernels **implicitly** map original vectors to higher dimensional spaces, take the dot product there, and hand the result back

- Linear kernel:
$$K(x, x') = x' \cdot x' = \sum_i x_i\, x'_i$$

- Quadratic kernel:
$$K(x, x') = (x \cdot x' + 1)^2$$

$$= \sum_{i,j} x_i x_j\, x'_i x'_j + 2 \sum_i x_i\, x'_i + 1$$

- RBF: infinite dimensional representation
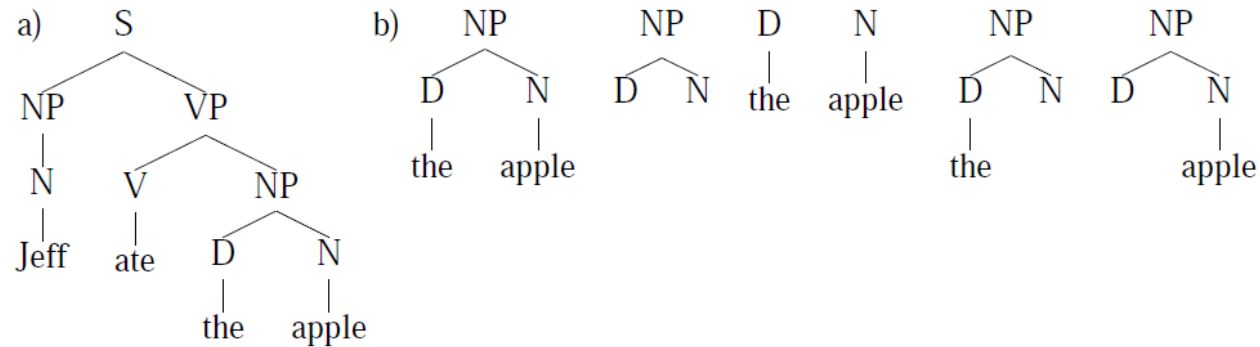$$K(x, x') = \exp(-||x - x'||^2)$$

- Discrete kernels: e.g. string kernels, tree kernels

- Want to compute number of common subtrees between T, T'

- Add up counts of all pairs of nodes n, n'
  - Base: if n, n' have different root productions, or are depth 0:

$$C(n_1, n_2) = 0$$

  - Base: if n, n' are share the same root production:

$$C(n_1, n_2) = \lambda \prod_{j=1}^{nc(n_1)} (1 + C(ch(n_1, j), ch(n_2, j)))$$

# Kernelized SVM (trust me)

**Primal formulation:**

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) - \xi_i$$

$$\mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})\left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right)$$

**Dual formulation:**

$$\min_{\alpha \geq 0} \quad \frac{1}{2}\left\|\sum_{i,y} \alpha_i(y)\left(f_i(y_i^*) - f_i(y)\right)\right\|_2^2 - \sum_{i,y}\alpha_i(y)\ell_i(y)$$

$$\forall i \quad \sum_y \alpha_i(y) = C$$

# Dual Formulation for SVMs

- We want to optimize: (separable case for now)

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})$$

- This is hard because of the constraints

- Solution: method of Lagrange multipliers

- The *Lagrangian* representation of this problem is:

$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \quad \Lambda(\mathbf{w}, \alpha) = \frac{1}{2}||\mathbf{w}||^2 - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}) \right)$$

- All we've done is express the constraints as an adversary which leaves our objective alone if we obey the constraints but ruins our objective if we violate any of them
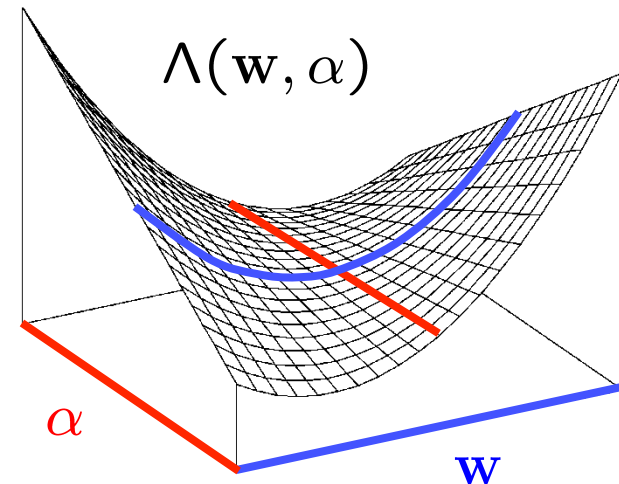
# Lagrange Duality

- We start out with a constrained optimization problem:

$$f(\mathbf{w}^*) \;=\; \min_{\mathbf{w}} \; f(\mathbf{w})$$

$$g(\mathbf{w}) \geq 0$$



$$\Lambda(\mathbf{w}, \alpha)$$

- We form the Lagrangian:

$$\Lambda(\mathbf{w}, \alpha) = f(\mathbf{w}) - \alpha \, g(\mathbf{w})$$

- This is useful because the constrained solution is a saddle point of $\Lambda$ (this is a general property):

$$f(\mathbf{w}^*) \;=\; \underbrace{\min_{\mathbf{w}} \max_{\alpha \geq 0} \Lambda(\mathbf{w}, \alpha)}_{\textit{Primal problem in } \mathbf{w}} = \underbrace{\max_{\alpha \geq 0} \min_{\mathbf{w}} \Lambda(\mathbf{w}, \alpha)}_{\textit{Dual problem in } \alpha}$$

# Dual Formulation II

- Duality tells us that

$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \quad \frac{1}{2}||\mathbf{w}||^2 - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}) \right)$$

has the same value as

$$Z(\alpha)$$

$$\max_{\alpha \geq 0} \min_{\mathbf{w}} \quad \frac{1}{2}||\mathbf{w}||^2 - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}) \right)$$

- This is useful because if we think of the $\alpha$'s as constants, we have an unconstrained min in w that we can solve analytically.
- Then we end up with an optimization over $\alpha$ instead of $\mathbf{w}$ (easier).

# Dual Formulation III

- Minimize the Lagrangian for fixed $\alpha$'s:

$$\Lambda(\mathbf{w}, \alpha) = \frac{1}{2}\|\mathbf{w}\|^2 - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y})\right)$$

$$\frac{\partial \Lambda(\mathbf{w}, \alpha)}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right)$$

$$\frac{\partial \Lambda(\mathbf{w}, \alpha)}{\partial \mathbf{w}} = 0 \quad \Longrightarrow \quad \mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right)$$

- So we have the Lagrangian as a function of only $\alpha$'s:

$$\min_{\alpha \geq 0} Z(\alpha) = \frac{1}{2} \left\| \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right) \right\|^2 - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \ell_i(\mathbf{y})$$

# Primal vs Dual SVM
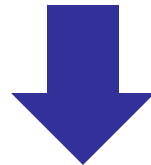
Primal formulation:

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) - \xi_i$$

$$\textcolor{red}{\mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) \left( \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y}) \right)}$$

Dual formulation:

$$\min_{\alpha \geq 0} \quad \frac{1}{2} \left\| \sum_{i,y} \alpha_i(y) \left( f_i(y_i^*) - f_i(y) \right) \right\|_2^2 - \sum_{i,y} \alpha_i(y) \ell_i(y)$$
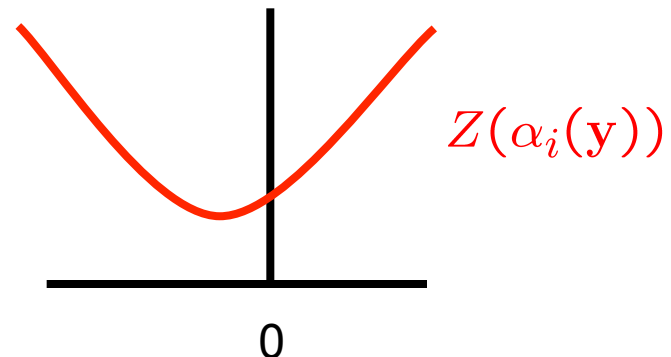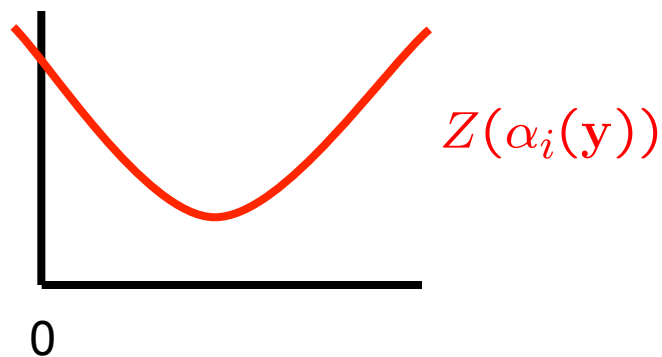
$$\forall i \quad \sum_y \alpha_i(y) = C$$

# Learning SVMs (Primal)

Primal formulation:

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) - \xi_i$$

$$\min_{w} \quad \frac{1}{2}\|w\|_2^2 + C\sum_i \left( \max_y \left( w^\top f_i(y) + \ell_i(y) \right) - w^\top f_i(y_i^*) \right)$$

# Learning SVMs (Primal)

Primal formulation:

$$\min_{w} \quad \frac{1}{2}\|w\|_2^2 + C \sum_i \left( \max_y \left( w^\top f_i(y) + \ell_i(y) \right) - w^\top f_i(y_i^*) \right)$$

Loss-augmented decode: $\quad \bar{y} = \operatorname{argmax}_y \left( w^\top f_i(y) + \ell_i(y) \right)$

$$\min_{w} \quad \frac{1}{2}\|w\|_2^2 + C \sum_i \left( w^\top f_i(\bar{y}) + \ell_i(\bar{y}) - w^\top f_i(y_i^*) \right)$$

$$\nabla_w = w + C \sum_i \left( f_i(\bar{y}) - f_i(y_i^*) \right)$$

Use general subgradient descent methods! (Adagrad)

# Learning SVMs (Dual)

- We want to find $\alpha$ which minimize

$$\min_{\alpha \geq 0} \quad \frac{1}{2} \left\| \sum_{i,y} \alpha_i(y)\big(f_i(y_i^*) - f_i(y)\big) \right\|_2^2 - \sum_{i,y} \alpha_i(y)\ell_i(y)$$

$$\forall i \quad \sum_y \alpha_i(y) = C$$

- This is a quadratic program:
  - Can be solved with general QP or convex optimizers
  - But they don't scale well to large problems
  - Cf. maxent models work fine with general optimizers (e.g. CG, L-BFGS)
- How would a special purpose optimizer work?

# Coordinate Descent I (Dual)

$$\min_{\alpha \geq 0} \quad \frac{1}{2} \left\| \sum_{i,y} \alpha_i(y)\big(f_i(y_i^*) - f_i(y)\big) \right\|_2^2 - \sum_{i,y} \alpha_i(y)\ell_i(y)$$

- Despite all the mess, $Z$ is just a quadratic in each $\alpha_i(\mathbf{y})$
- Coordinate descent: optimize one variable at a time

$Z(\alpha_i(\mathbf{y}))$

$Z(\alpha_i(\mathbf{y}))$

0

0

- If the unconstrained argmin on a coordinate is negative, just clip to zero...

# Coordinate Descent II (Dual)

- Ordinarily, treating coordinates independently is a bad idea, but here the update is very fast and simple

$$\alpha_i(\mathbf{y}) \leftarrow \max\left(0, \alpha_i(\mathbf{y}) + \frac{\ell_i(\mathbf{y}) - \mathbf{w}^\top\left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right)}{\left\|\left(\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})\right)\right\|^2}\right)$$

- So we visit each axis many times, but each visit is quick

- This approach works fine for the separable case
- For the non-separable case, we just gain a simplex constraint and so we need slightly more complex methods (SMO, exponentiated gradient)

$$\forall i, \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C$$

# What are the Alphas?

- Each candidate corresponds to a primal constraint

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) - \xi_i$$

*Support vectors*

- In the solution, an $\alpha_i(\mathbf{y})$ will be:
  - Zero if that constraint is inactive
  - Positive if that constrain is active
  - i.e. positive on the support vectors

- Support vectors contribute to weights:

$$\mathbf{w} = \sum_{i, \mathbf{y}} \alpha_i(\mathbf{y}) \left( \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y}) \right)$$

# Structure

# Handwriting recognition

**x**                                                  **y**

 → **brace**

Sequential structure

**x**

**y**

*The screen was
a sea of red*

➡️



Recursive structure

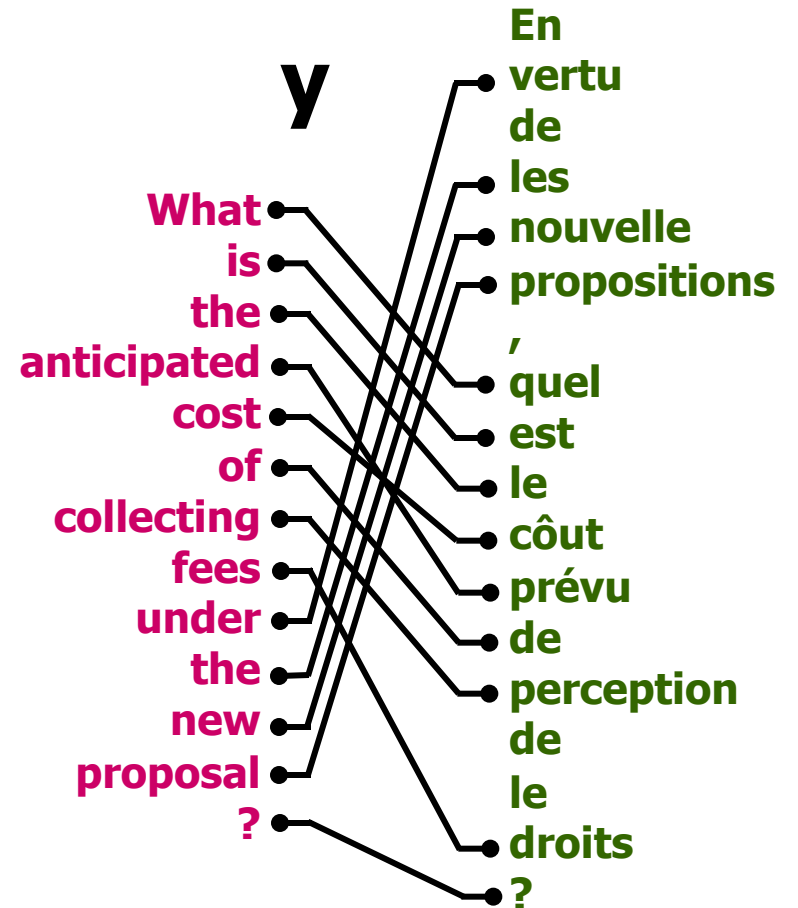# Bilingual Word Alignment

**x**

**What is the anticipated cost of collecting fees under the new proposal?**

**En vertu de nouvelle propositions, quel est le côut prévu de perception de les droits?**

**y**



Combinatorial structure

# Structured Models

$$prediction(\mathbf{x}, \mathbf{w}) = \arg\max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} score(\mathbf{y}, \mathbf{w})$$

space of feasible outputs

Assumption:

$$score(\mathbf{y}, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}(\mathbf{y}) = \sum_p \mathbf{w}^\top \mathbf{f}(\mathbf{y}_p)$$

Score is a sum of local "part" scores

Parts = nodes, edges, productions

$$\sum_{y_{jk} \in \mathbf{y}} \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}_{jk}) = \mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y})$$

What
is
the
anticipated
cost
of
collecting
fees
under
the
new
proposal
?

En
vertu
de
les
nouvelle
propositions
,
quel
est
le
côut
prévu
de
perception
de
le
droits
?

j

k

$y_{jk}$

$\mathbf{f}(\mathbf{x}_{jk})$

- association

- position

- orthography

# Efficient Decoding

- Common case: you have a black box which computes

$$\text{prediction}(\mathbf{x}) = \underset{\mathbf{y} \in \mathcal{Y}(\mathbf{x})}{\arg\max} \mathbf{w}^\top \mathbf{f}(\mathbf{y})$$

  at least approximately, and you want to learn w

- Easiest option is the structured perceptron [Collins 01]
  - Structure enters here in that the search for the best y is typically a combinatorial algorithm (dynamic programming, matchings, ILPs, A*…)
  - Prediction is structured, learning update is not

Remember our primal margin objective?

$$\min_{w} \quad \frac{1}{2}\|w\|_2^2 + C\sum_i \left( \max_y \left( w^\top f_i(y) + \ell_i(y) \right) - w^\top f_i(y_i^*) \right)$$

Still applies with structured output space!

# Structured Margin (Primal)

Just need efficient loss-augmented decode:

$$\bar{y} = \text{argmax}_y \left( w^\top f_i(y) + \ell_i(y) \right)$$

$$\min_w \quad \frac{1}{2} \|w\|_2^2 + C \sum_i \left( w^\top f_i(\bar{y}) + \ell_i(\bar{y}) - w^\top f_i(y_i^*) \right)$$

$$\nabla_w = w + C \sum_i \left( f_i(\bar{y}) - f_i(y_i^*) \right)$$

Still use general subgradient descent methods! (Adagrad)

# Structured Margin (Dual)

- Remember the constrained version of primal:

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

$$\forall i, \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) - \xi_i$$

- Dual has a variable for every constraint here

- We want:

$$\arg\max_{\mathbf{y}} \ \mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \mathbf{y}) \ = \ \text{``brace''}$$

- Equivalently:

$$\mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``brace''}) \ > \ \mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``aaaaa''})$$

$$\mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``brace''}) \ > \ \mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``aaaab''})$$

$$\ldots$$

$$\mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``brace''}) \ > \ \mathbf{w}^\top \mathbf{f}(\boxed{\textit{brace}}, \text{``zzzzz''})$$

**a lot!**

# Parsing example

- We want:

$$\arg\max_y \ w^\top f(\ \text{`It was red'}\ , y) \ = \ \text{[tree: S → A B, C D]}$$

- Equivalently:

$$w^\top f(\text{`It was red'}, \text{[tree]}) \ > \ w^\top f(\text{`It was red'}, \text{[tree]})$$

$$w^\top f(\text{`It was red'}, \text{[tree]}) \ > \ w^\top f(\text{`It was red'}, \text{[tree]})$$

$$\ldots$$

$$w^\top f(\text{`It was red'}, \text{[tree]}) \ > \ w^\top f(\text{`It was red'}, \text{[tree]})$$

**a lot!**

# Alignment example

- We want:

$$\arg\max_y \ \mathbf{w}^\top \mathbf{f}(\ \substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, y) \ = \ \substack{1 \leftrightarrow 1\\ 2 \leftrightarrow 2\\ 3 \leftrightarrow 3}$$

- Equivalently:

$$\mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \leftrightarrow 1\\ 2 \leftrightarrow 2\\ 3 \leftrightarrow 3}) \ > \ \mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \leftrightarrow 1\\ 2 \times 2\\ 3 \leftrightarrow 3})$$

$$\mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \leftrightarrow 1\\ 2 \leftrightarrow 2\\ 3 \leftrightarrow 3}) \ > \ \mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \times 1\\ 2 \times 2\\ 3 \leftrightarrow 3})$$

$$\dots$$

$$\mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \leftrightarrow 1\\ 2 \leftrightarrow 2\\ 3 \leftrightarrow 3}) \ > \ \mathbf{w}^\top \mathbf{f}(\substack{\text{‘What is the’}\\ \text{‘Quel est le’}}, \substack{1 \times 1\\ 2 \times 2\\ 3 \times 3})$$

**a lot!**

# Cutting Plane (Dual)

- A constraint induction method [Joachims et al 09]
  - Exploits that the number of constraints you actually need per instance is typically very small
  - Requires (loss-augmented) primal-decode only

- Repeat:
  - Find the most violated constraint for an instance:

$$\forall \mathbf{y} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})$$
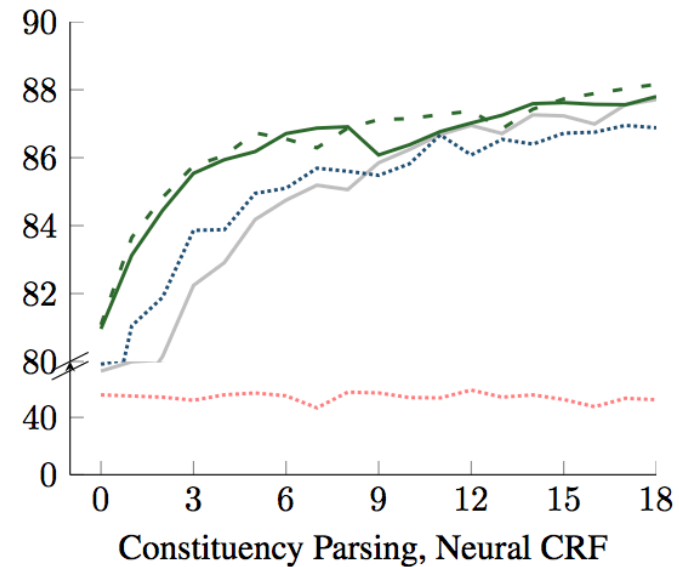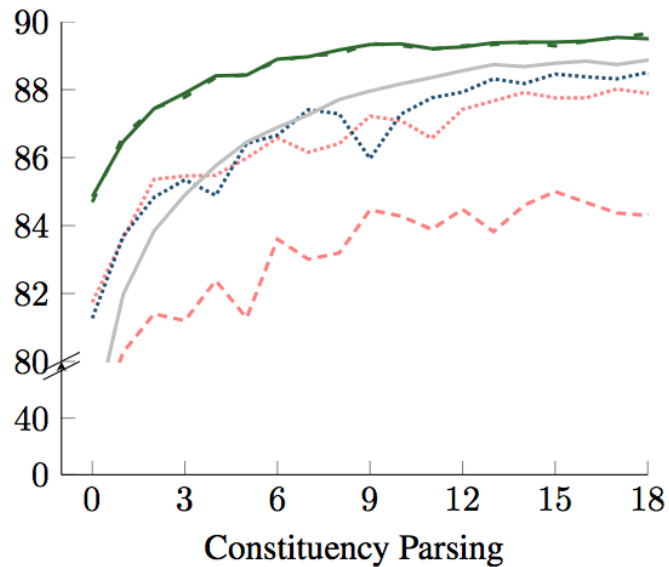
$$\arg\max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})$$

  - Add this constraint and resolve the (non-structured) QP (e.g. with SMO or other QP solver)

# Comparison

| 8 | | | |
|---|---|---|---|
| | Oct 20 | Structured Classification III | |
| 9 | Oct 25 | Structured Classification IV | J+M 16, 18, 19, Adagrad, Subgradient SVM |



Constituency Parsing

Constituency Parsing, Neural CRF

| Margin | - - - Cutting Plane |
|---|---|
| | ······· Online Cutting Plane |
| | - - - Online Primal Subgradient & $L_1$ |
| | ——— Online Primal Subgradient & $L_2$ |
| Mistake Driven | - - - Averaged Perceptron |
| | ······· MIRA |
| | - - Averaged MIRA (MST built-in) |
| Llhood | ——— Stochastic Gradient Descent |

# Option 0: Reranking

Input

N-Best List
(e.g. n=100)

Output

x =
"The screen was a sea of red."

Baseline Parser

Non-Structured Classification

- **Advantages:**
  - Directly reduce to non-structured case
  - No locality restriction on features

$$\mathbf{f}\left( \begin{array}{c} \text{[parse tree]} \end{array} \right) =$$

- **Disadvantages:**
  - Stuck with errors of baseline parser
  - Baseline system must produce n-best lists
  - But, feedback is possible [McCloskey, Charniak, Johnson 2006]