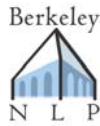


# Natural Language Processing



Classification II  
Dan Klein – UC Berkeley

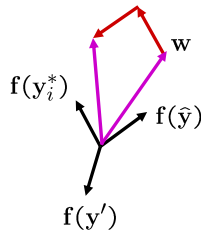
# Classification



## Linear Models: Perceptron

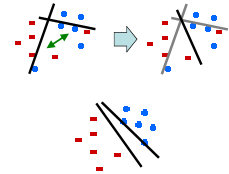
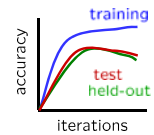
- The perceptron algorithm
  - Iteratively processes the training set, reacting to training errors
  - Can be thought of as trying to drive down training error
- The (online) perceptron algorithm:
  - Start with zero weights  $w$
  - Visit training instances one by one
    - Try to classify
 
$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} w^T f(y)$$
    - If correct, no change!
    - If wrong: adjust weights
 
$$w \leftarrow w + f(y_i^*)$$

$$w \leftarrow w - f(\hat{y})$$



## Issues with Perceptrons

- Overtraining: test / held-out accuracy usually rises, then falls
  - Overtraining isn't the typically discussed source of overfitting, but it can be important
- Regularization: if the data isn't separable, weights often thrash around
  - Averaging weight vectors over time can help (averaged perceptron)
  - [Freund & Schapire 99, Collins 02]
- Mediocre generalization: finds a "barely" separating solution

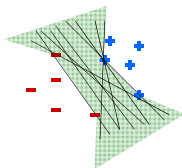


## Problems with Perceptrons

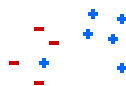
- Perceptron "goal": separate the training data

$$\forall i, \forall y \neq y^i \quad w^T f_i(y^i) \geq w^T f_i(y)$$

1. This may be an entire feasible space



2. Or it may be impossible



# Margin

## Objective Functions

- What do we want from our weights?
  - Depends!
  - So far: minimize (training) errors:

$$\sum_i \text{step} \left( w^\top f_i(y_i^*) - \max_{y \neq y_i^*} w^\top f_i(y) \right)$$

- This is the “zero-one loss”
  - Discontinuous, minimizing is NP-complete
  - Not really what we want anyway
- Maximum entropy and SVMs have other objectives related to zero-one loss

## Linear Separators

- Which of these linear separators is optimal?

## Classification Margin (Binary)

- Distance of  $x_i$  to separator is its margin,  $m_i$
- Examples closest to the hyperplane are **support vectors**
- Margin  $\gamma$**  of the separator is the minimum  $m$

## Classification Margin

- For each example  $x_i$  and possible mistaken candidate  $y$ , we avoid that mistake by a margin  $m_i(y)$  (with zero-one loss)
 
$$m_i(y) = w^\top f_i(y_i^*) - w^\top f_i(y)$$
- Margin  $\gamma$  of the entire separator is the minimum  $m$ 

$$\gamma = \min_i \left( w^\top f_i(y_i^*) - \max_{y \neq y_i^*} w^\top f_i(y) \right)$$
- It is also the largest  $\gamma$  for which the following constraints hold
 
$$\forall i, \forall y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + \gamma \ell_i(y)$$

## Maximum Margin

- Separable SVMs: find the max-margin  $w$

$$\max_{\|w\|=1} \gamma \quad \ell_i(y) = \begin{cases} 0 & \text{if } y = y_i^* \\ 1 & \text{if } y \neq y_i^* \end{cases}$$

$$\forall i, \forall y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + \gamma \ell_i(y)$$

- Can stick this into Matlab and (slowly) get an SVM
- Won't work (well) if non-separable

## Why Max Margin?

- Why do this? Various arguments:
  - Solution depends only on the boundary cases, or **support vectors** (but remember how this diagram is broken!)
  - Solution robust to movement of support vectors
  - Sparse solutions (features not in support vectors get zero weight)
  - Generalization bound arguments
  - Works well in practice for many problems**

Support vectors

## Max Margin / Small Norm

- Reformulation: find the smallest  $w$  which separates data

Remember this condition?  $\xrightarrow{\max_{\|w\|=1} \gamma}$

$$\forall i, y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + \gamma \ell_i(y)$$

- $\gamma$  scales linearly in  $w$ , so if  $\|w\|$  isn't constrained, we can take any separating  $w$  and scale up our margin

$$\gamma = \min_{i, y \neq y_i} [w^\top f_i(y_i^*) - w^\top f_i(y)] / \ell_i(y)$$

- Instead of fixing the scale of  $w$ , we can fix  $\gamma = 1$

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall i, y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + 1 \ell_i(y)$$

## Soft Margin Classification

- What if the training set is not linearly separable?
- Slack variables  $\xi_i$  can be added to allow misclassification of difficult or noisy examples, resulting in a *soft margin* classifier

## Maximum Margin

Note: exist other choices of how to penalize slacks!

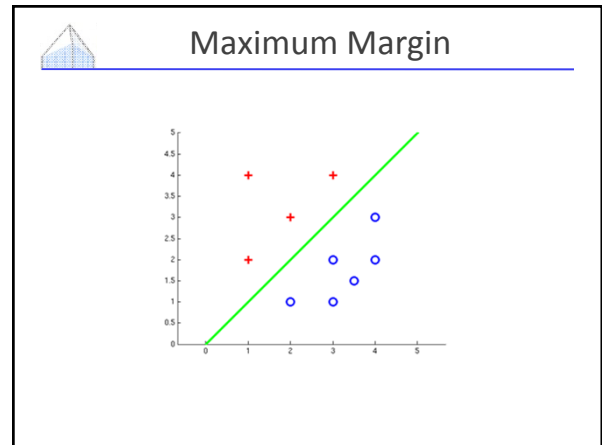
- Non-separable SVMs
  - Add slack to the constraints
  - Make objective pay (linearly) for slack:

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

$$\forall i, y, \quad w^\top f_i(y_i^*) + \xi_i \geq w^\top f_i(y) + \ell_i(y)$$

- $C$  is called the *capacity* of the SVM – the smoothing knob

- Learning:
  - Can still stick this into Matlab if you want
  - Constrained optimization is hard; better methods!
  - We'll come back to this later



## Likelihood

## Linear Models: Maximum Entropy

- Maximum entropy (logistic regression)
  - Use the scores as probabilities:
 
$$P(y|x, w) = \frac{\exp(w^\top f(y))}{\sum_{y'} \exp(w^\top f(y'))}$$

← Make Nonlinearize
  - Maximize the (log) conditional likelihood of training data

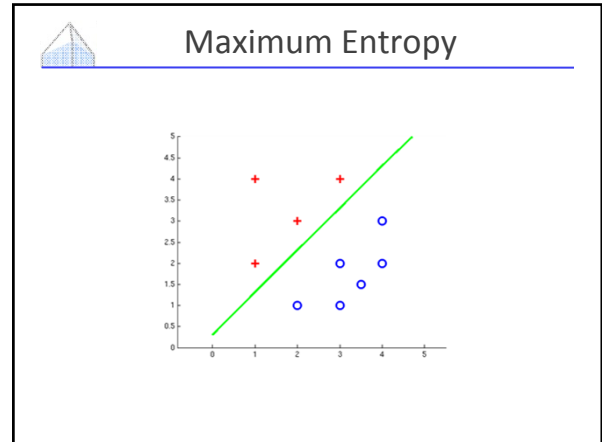
$$L(w) = \log \prod_i P(y_i^* | x_i, w) = \sum_i \log \left( \frac{\exp(w^\top f_i(y_i^*))}{\sum_y \exp(w^\top f_i(y))} \right)$$

$$= \sum_i \left( w^\top f_i(y_i^*) - \log \sum_y \exp(w^\top f_i(y)) \right)$$

## Maximum Entropy II

- Motivation for maximum entropy:
  - Connection to maximum entropy principle (sort of)
  - Might want to do a good job of being uncertain on noisy cases...
  - ... in practice, though, posteriors are pretty peaked
- Regularization (smoothing)
 
$$\max_{\mathbf{w}} \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(y_i^*) - \log \sum_y \exp(\mathbf{w}^\top \mathbf{f}_i(y)) \right) - k \|\mathbf{w}\|^2$$

$$\min_{\mathbf{w}} k \|\mathbf{w}\|^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(y_i^*) - \log \sum_y \exp(\mathbf{w}^\top \mathbf{f}_i(y)) \right)$$



## Loss Comparison

## Log-Loss

- If we view maxent as a minimization problem:
 
$$\min_{\mathbf{w}} k \|\mathbf{w}\|^2 + \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(y_i^*) - \log \sum_y \exp(\mathbf{w}^\top \mathbf{f}_i(y)) \right)$$
- This minimizes the "log loss" on each example
 
$$-\left( \mathbf{w}^\top \mathbf{f}_i(y_i^*) - \log \sum_y \exp(\mathbf{w}^\top \mathbf{f}_i(y)) \right) = -\log P(y_i^* | \mathbf{x}_i, \mathbf{w})$$

step  $\left( \mathbf{w}^\top \mathbf{f}_i(y_i^*) - \max_{y \neq y_i^*} \mathbf{w}^\top \mathbf{f}_i(y) \right)$
- One view: log loss is an *upper bound* on zero-one loss

## Remember SVMs...

- We had a **constrained** minimization
 
$$\min_{\mathbf{w}, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

$$\forall i, y, \quad \mathbf{w}^\top \mathbf{f}_i(y_i^*) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(y) + \xi_i(y)$$
- ...but we can solve for  $\xi_i$ 

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

$$\forall i, \quad \xi_i = \max(\mathbf{w}^\top \mathbf{f}_i(y) + \xi_i(y) - \mathbf{w}^\top \mathbf{f}_i(y_i^*))$$
- Giving
 
$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \left( \max_y (\mathbf{w}^\top \mathbf{f}_i(y) + \xi_i(y)) - \mathbf{w}^\top \mathbf{f}_i(y_i^*) \right)$$

## Hinge Loss

Plot really only right in binary case

- Consider the per-instance objective:
 
$$\min_{\mathbf{w}} k \|\mathbf{w}\|^2 + \sum_i \left( \max_y (\mathbf{w}^\top \mathbf{f}_i(y) + \xi_i(y)) - \mathbf{w}^\top \mathbf{f}_i(y_i^*) \right)$$
- This is called the "**hinge loss**"
  - Unlike **maxent / log loss**, you stop gaining objective once the true label wins by enough
  - You can start from here and derive the SVM objective
  - Can solve directly with sub-gradient decent (e.g. Pegasos: Shalev-Shwartz et al 07)

$\mathbf{w}^\top \mathbf{f}_i(y_i^*) - \max_{y \neq y_i^*} (\mathbf{w}^\top \mathbf{f}_i(y))$

### Max vs "Soft-Max" Margin

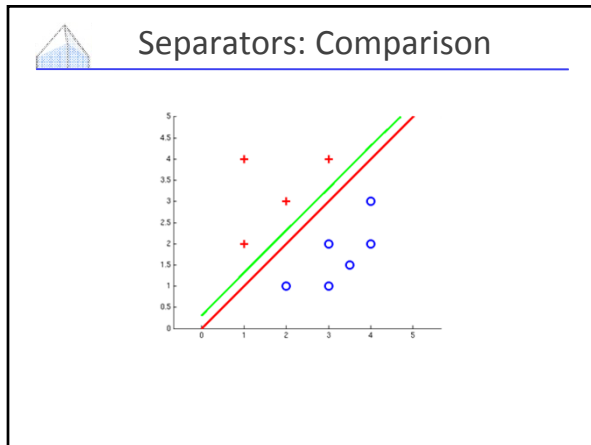
- SVMs:
 
$$\min_w k \|w\|^2 - \sum_i \left( w^T f_i(y_i^*) - \max_{y \neq y_i^*} (w^T f_i(y) + \epsilon_i(y)) \right)$$

You can make this zero
- Maxent:
 
$$\min_w k \|w\|^2 - \sum_i \left( w^T f_i(y_i^*) - \log \sum_y \exp(w^T f_i(y)) \right)$$

... but not this one
- Very similar! Both try to make the true score better than a function of the other scores
  - The SVM tries to beat the augmented runner-up
  - The Maxent classifier tries to beat the "soft-max"

### Loss Functions: Comparison

- Zero-One Loss
 
$$\sum_i \text{step} \left( w^T f_i(y_i^*) - \max_{y \neq y_i^*} w^T f_i(y) \right)$$
- Hinge
 
$$\sum_i \left( w^T f_i(y_i^*) - \max_{y \neq y_i^*} (w^T f_i(y) + \epsilon_i(y)) \right)$$
- Log
 
$$\sum_i \left( w^T f_i(y_i^*) - \log \sum_y \exp(w^T f_i(y)) \right)$$



### Conditional vs Joint Likelihood

### Example: Sensors

**Reality**

Raining: Sunny:

$P(+,+,r) = 3/8$      $P(-,+,r) = 1/8$      $P(+,+,s) = 1/8$      $P(-,+,s) = 3/8$

**NB Model**

Raining?

M1    M2

**NB FACTORS:**

- $P(s) = 1/2$
- $P(+|s) = 1/4$
- $P(+|r) = 3/4$

**PREDICTIONS:**

- $P(r,+,+) = (1/2)(3/4)(3/4)$
- $P(s,+,+) = (1/2)(1/4)(1/4)$
- $P(r|+,+) = 9/10$
- $P(s|+,+) = 1/10$

### Example: Stoplights

**Reality**

Lights Working: Lights Broken:

$P(g,r,w) = 3/7$      $P(r,g,w) = 3/7$      $P(r,r,b) = 1/7$

**NB Model**

Working?

NS    EW

**NB FACTORS:**

- $P(w) = 6/7$
- $P(r|w) = 1/2$
- $P(g|w) = 1/2$
- $P(b) = 1/7$
- $P(r|b) = 1$
- $P(g|b) = 0$



## Example: Stoplights

- What does the model say when both lights are red?
  - $P(b,r,r) = (1/7)(1)(1) = 1/7 = 4/28$
  - $P(w,r,r) = (6/7)(1/2)(1/2) = 6/28 = 6/28$
  - $P(w|r,r) = 6/10!$
- We'll guess that  $(r,r)$  indicates lights are working!
- Imagine if  $P(b)$  were boosted higher, to  $1/2$ :
  - $P(b,r,r) = (1/2)(1)(1) = 1/2 = 4/8$
  - $P(w,r,r) = (1/2)(1/2)(1/2) = 1/8 = 1/8$
  - $P(w|r,r) = 1/5!$
- Changing the parameters bought accuracy at the expense of data likelihood

## Duals and Kernels



## Nearest-Neighbor Classification

- Nearest neighbor, e.g. for digits:
  - Take new example
  - Compare to all training examples
  - Assign based on closest example

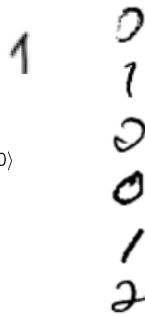
- Encoding: image is vector of intensities:

$$1 = \langle 0.0 \ 0.0 \ 0.3 \ 0.8 \ 0.7 \ 0.1 \dots 0.0 \rangle$$

- Similarity function:

- E.g. dot product of two images' vectors

$$\text{sim}(x, y) = x^T y = \sum_i x_i y_i$$



## Non-Parametric Classification

- Non-parametric: more examples means (potentially) more complex classifiers
- How about K-Nearest Neighbor?
  - We can be a little more sophisticated, averaging several neighbors
  - But, it's still not really error-driven learning
  - The magic is in the distance function
- Overall: we can exploit rich similarity functions, but not objective-driven learning



## A Tale of Two Approaches...

- Nearest neighbor-like approaches
  - Work with data through similarity functions
  - No explicit "learning"
- Linear approaches
  - Explicit training to reduce empirical error
  - Represent data through features
- Kernelized linear models
  - Explicit training, but driven by similarity!
  - Flexible, powerful, very very slow



## The Perceptron, Again

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

$$\hat{y} = \arg \max_{y \in \mathcal{Y}(x)} w^T f_i(y)$$

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

$$w \leftarrow w + f_i(y_i^*)$$

$$w \leftarrow w - f_i(\hat{y})$$

$$w \leftarrow w + (f_i(y_i^*) - f_i(\hat{y}))$$

$$w \leftarrow w + \Delta_i(\hat{y}) \quad \text{mistake vectors}$$

## Perceptron Weights

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

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

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

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

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

$$\alpha_i = \langle \alpha_i(\mathbf{y}_1) \ \alpha_i(\mathbf{y}_2) \ \dots \ \alpha_i(\mathbf{y}_n) \rangle$$

## Dual Perceptron

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

- Track mistake counts rather than weights
- Start with zero counts ( $\alpha$ )
- For each instance  $\mathbf{x}$ 
  - Try to classify

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

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

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

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

## Dual / Kernelized Perceptron

- How to classify an example  $\mathbf{x}$ ?

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

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

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

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

- 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

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

- 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

## Tree Kernels [Collins and Duffy 01]

a)

b)

- 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' share the same root production:
 
$$C(n_1, n_2) = \lambda \prod_{j=1}^{n_c(n_1)} (1 + C(ch(n_1, j), ch(n_2, j)))$$

44

## Dual Formulation for SVMs

- We want to optimize: (separable case for now)
 
$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2$$

$$\forall i, y \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})$$
- This is hard because of the constraints
- Solution: method of Lagrange multipliers
- The Lagrangian representation of this problem is:
 
$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \Lambda(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}))$$
- 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:
 
$$f(\mathbf{w}^*) = \min_{\mathbf{w}} f(\mathbf{w})$$

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

## Dual Formulation II

- Duality tells us that
 
$$\min_{\mathbf{w}} \max_{\alpha \geq 0} \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}))$$
- has the same value as
 
$$\max_{\alpha \geq 0} \underbrace{\min_{\mathbf{w}} \left( \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y})) \right)}_{Z(\alpha)}$$
- This is useful because if we think of the  $\alpha$ 's as constants, we have an unconstrained min in  $\mathbf{w}$  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:
 
$$\Lambda(\mathbf{w}, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) - \ell_i(\mathbf{y}))$$

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

$$\frac{\partial \Lambda(\mathbf{w}, \alpha)}{\partial \mathbf{w}} = 0 \Rightarrow \mathbf{w} = \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y}))$$
- So we have the Lagrangian as a function of only  $\alpha$ 's:
 
$$\min_{\alpha \geq 0} Z(\alpha) = \frac{1}{2} \left\| \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})) \right\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) \ell_i(\mathbf{y})$$

## Back to Learning SVMs

- We want to find  $\alpha$  which minimize
 
$$\min_{\alpha \geq 0} \Lambda(\alpha) = \frac{1}{2} \left\| \sum_{i, y} \alpha_i(\mathbf{y}) (\mathbf{f}_i(\mathbf{y}_i^*) - \mathbf{f}_i(\mathbf{y})) \right\|^2 - \sum_{i, y} \alpha_i(\mathbf{y}) \ell_i(\mathbf{y})$$

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

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

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

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

## Coordinate Descent II

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

$$\alpha_i(y) \leftarrow \max \left( 0, \alpha_i(y) + \frac{\ell_i(y) - w^T (f_i(y_i^*) - f_i(y))}{\|f_i(y_i^*) - f_i(y)\|^2} \right)$$

- 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)

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

## What are the Alphas?

- Each candidate corresponds to a primal constraint

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

$$\forall i, y \quad w^T f_i(y_i^*) \geq w^T f_i(y) + \ell_i(y) - \xi_i$$

- In the solution, an  $\alpha_i(y)$  will be:
  - Zero if that constraint is inactive
  - Positive if that constraint is active
  - i.e. positive on the support vectors
- Support vectors contribute to weights:

$$w = \sum_{i,y} \alpha_i(y) (f_i(y_i^*) - f_i(y))$$

Structure

## Handwriting recognition

**x** **y**

Sequential structure

[Slides: Taskar and Klein 05]

## CFG Parsing

**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 coût prévu de perception de les droits?

**y**

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

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

Combinatorial structure

### Structured Models

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

space of feasible outputs

Assumption:

$$score(y, w) = w^\top f(y) = \sum_p w^\top f(y_p)$$

Score is a sum of local "part" scores  
Parts = nodes, edges, productions

### CFG Parsing

$$P(y | x) \propto \prod_{A \rightarrow \alpha \in (x, y)} \phi(A \rightarrow \alpha)$$

#(NP → DT NN) ...  
#(PP → IN NP) ...  
#(NN → 'sea')

$f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$

$$\prod_{A \rightarrow \alpha \in (x, y)} \exp\{w^\top f(A \rightarrow \alpha)\} = \exp\{w^\top f(x, y)\}$$

### Bilingual word alignment

$$\sum_{y_{jk} \in \mathcal{Y}} w^\top f(x_{jk}) = w^\top f(x, y)$$

$f(x_{jk})$

- association
- position
- orthography

### Option 0: Reranking

[e.g. Charniak and Johnson 05]

Input: "The screen was a sea of red."

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

Output

Baseline Parser → Non-Structured Classification

60

### Reranking

- Advantages:
  - Directly reduce to non-structured case
  - No locality restriction on features
- Disadvantages:
  - Stuck with errors of baseline parser
  - Baseline system must produce n-best lists
  - But, feedback is possible [McCloskey, Charniak, Johnson 2006]

$$f(\text{Parse Tree}) =$$



## Efficient Primal Decoding

- Common case: you have a black box which computes

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

at least approximately, and you want to learn  $w$

- Many learning methods require more (expectations, dual representations, k-best lists), but the most commonly used options do not
- 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

62



## Structured Margin

- Remember the margin objective:

$$\min_w \frac{1}{2} \|w\|^2$$

$$\forall_i, y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + \ell_i(y)$$

- This is still defined, but lots of constraints



## Full Margin: OCR

- We want:

$$\arg \max_y w^\top f(\text{brace}, y) = \text{"brace"}$$

- Equivalently:

$$\begin{aligned} w^\top f(\text{brace}, \text{"brace"}) &> w^\top f(\text{brace}, \text{"aaaa"}) \\ w^\top f(\text{brace}, \text{"brace"}) &> w^\top f(\text{brace}, \text{"aaaab"}) \\ &\dots \\ w^\top f(\text{brace}, \text{"brace"}) &> w^\top f(\text{brace}, \text{"zzzz"}) \end{aligned} \quad \left. \vphantom{\begin{aligned} & \\ & \\ & \\ & \end{aligned}} \right\} \text{a lot!}$$



## Parsing example

- We want:

$$\arg \max_y w^\top f(\text{"It was red"}, y) = \hat{A}_{c_b}^{\hat{A}_b}$$

- Equivalently:

$$\begin{aligned} w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_b}) &> w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_a}) \\ w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_b}) &> w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_c}) \\ &\dots \\ w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_b}) &> w^\top f(\text{"It was red"}, \hat{A}_{c_b}^{\hat{A}_d}) \end{aligned} \quad \left. \vphantom{\begin{aligned} & \\ & \\ & \\ & \end{aligned}} \right\} \text{a lot!}$$



## Alignment example

- We want:

$$\arg \max_y w^\top f(\text{"What is the"}, \text{"Quel est le'}) = \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}$$

- Equivalently:

$$\begin{aligned} w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) &> w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) \\ w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) &> w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) \\ &\dots \\ w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) &> w^\top f(\text{"What is the"}, \begin{matrix} 1 \leftrightarrow 1 \\ 2 \leftrightarrow 2 \\ 3 \leftrightarrow 3 \end{matrix}) \end{aligned} \quad \left. \vphantom{\begin{aligned} & \\ & \\ & \\ & \end{aligned}} \right\} \text{a lot!}$$



## Cutting Plane

- 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:

$$\forall y \quad w^\top f_i(y_i^*) \geq w^\top f_i(y) + \ell_i(y)$$

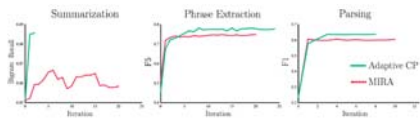
$$\arg \max_y w^\top f_i(y) + \ell_i(y)$$

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



## Cutting Plane

- Some issues:
  - Can easily spend too much time solving QPs
  - Doesn't exploit shared constraint structure
  - In practice, works pretty well; fast like MIRA, more stable, no averaging



## M3Ns

- Another option: express all constraints in a packed form
  - Maximum margin Markov networks [Taskar et al 03]
  - Integrates solution structure deeply into the problem structure
- Steps
  - Express inference over constraints as an LP
  - Use duality to transform minimax formulation into min-min
  - Constraints factor in the dual along the same structure as the primal; alphas essentially act as a dual "distribution"
  - Various optimization possibilities in the dual



## Likelihood, Structured

$$L(\mathbf{w}) = -k\|\mathbf{w}\|^2 + \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}_i^*) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -2k\mathbf{w} + \sum_i \left( \mathbf{f}_i(\mathbf{y}_i^*) - \sum_{\mathbf{y}} P(\mathbf{y}|\mathbf{x}_i) \mathbf{f}_i(\mathbf{y}) \right)$$

- Structure needed to compute:
  - Log-normalizer
  - Expected feature counts
    - E.g. if a feature is an indicator of DT-NN then we need to compute posterior marginals  $P(\text{DT-NN}|\text{sentence})$  for each position and sum
- Also works with latent variables (more later)