# Max-Margin Methods for NLP:
## Estimation, Structure, and Applications

Dan Klein and Ben Taskar
EECS Department
UC Berkeley

---

# Introduction

- Much of NLP can be seen as making decisions
  - About structured analyses (sequences, trees, graphs)
  - On the basis of multiple information sources, or *features* (words, word classes, tree configurations, etc.)

- Widespread adoption of discriminative methods
  - Use of arbitrary features
  - Various formulations: maxent, SVM, perceptron
  - Common use: local discriminative decisions, possibly chained
  - Relatively new: global methods which exploit model structure (CRFs, max-margin networks)

- This tutorial will cover:
  - Part I: Flat max-margin methods (SVMs)
  - Part II: Structured max-margin methods (sequences, trees, matchings)

# Outline

- ## Part I: Flat Classification
  - Linear classifiers and loss functions
  - Primal and dual SVM formulations
  - Training SVMs

- ## Part II: Structured Classification
  - Structured linear classifiers
  - Factored learning formulations
  - Experimental results

---

# Example: Text Classification

- We want to classify documents into categories

| DOCUMENT | CATEGORY |
| --- | --- |
| … win the election … | *POLITICS* |
| … win the game … | *SPORTS* |
| … see a movie … | *OTHER* |

- Classically, do this on the basis of words in the document, but other information sources are potentially relevant:
  - Document length
  - Average word length
  - Document's source
  - Burstiness of new words in document
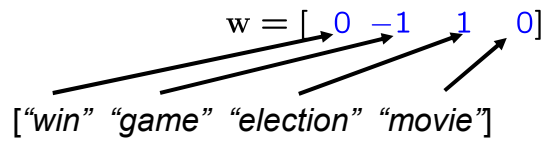
# Some Definitions

| | | |
|---|---|---|
| INPUTS | $\mathbf{x}^i$ | *… win the election …* |
| TRUE OUTPUTS | $\mathbf{y}^i$ | *POLITICS* |
| OUTPUT SPACE | $\mathcal{Y}$ | *SPORTS, POLITICS, OTHER* |
| ANY OUTPUTS | $\mathbf{y}$ | *SPORTS, POLITICS, OTHER* |

# Binary Linear Models

- Two Classes   *POLITICS* = +, *SPORTS* = -
- Features   $\mathbf{f}(\text{…win the election…}) = [\ \ 1 \quad 0 \quad 1 \quad 0]$
- Weights   $\mathbf{w} = [\ \ 0 \quad -1 \quad 1 \quad 0]$

  [*"win" "game" "election" "movie"*]

- Prediction rule

$$\text{prediction}(\mathbf{x}, \mathbf{w}) =$$
$$\begin{cases} + & \text{if } \mathbf{w}^\top \mathbf{f}(\mathbf{x}) \geq 0 \\ - & \text{if } \mathbf{w}^\top \mathbf{f}(\mathbf{x}) < 0 \end{cases}$$

$\mathbf{w}^\top \mathbf{f}$

$\mathbf{w}^\top \mathbf{f} > 0$

$\mathbf{w}^\top \mathbf{f} < 0$   $\mathbf{w}^\top \mathbf{f} = 0$

# Multiclass Linear Models

- **Multiple Classes** *SPORTS, POLITICS, OTHER*

$\mathbf{f}_i(POLITICS) =$ [ 0  0  0  0  1  0  1  0  0  0  0  0]
$\mathbf{f}_i(SPORTS) =$ [ 1  0  1  0  0  0  0  0  0  0  0  0]
$\mathbf{f}_i(OTHER) =$ [ 0  0  0  0  0  0  0  0  1  0  1  0]
$\mathbf{w} =$ [ 1  1  −1  −2  1  −1  1  −2  −2  −1  −1  1]

[*"win"*∧SPORTS  *"game"*∧SPORTS  *"election"*∧SPORTS  *"movie"*∧SPORTS ]

$$\mathbf{f(x,y)} = [\ 0 \quad 0 \quad \cdots \quad \mathbf{f(x)} \quad \cdots \quad 0\ ]$$
$$\mathbf{w} = [\ \mathbf{w_0} \quad \mathbf{w_1} \quad \cdots \quad \mathbf{w_y} \quad \cdots \quad \mathbf{w_k}\ ]$$

---

# Multiclass Linear Models

$$\mathbf{f(x,y)} = [\ 0 \quad 0 \quad \cdots \quad \mathbf{f(x)} \quad \cdots \quad 0\ ]$$
$$\mathbf{w} = [\ \mathbf{w_0} \quad \mathbf{w_1} \quad \cdots \quad \mathbf{w_y} \quad \cdots \quad \mathbf{w_k}\ ]$$

- **Scores and Predictions**

$$score(\mathbf{x}^i, \mathbf{y}, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) = \mathbf{w_y}^\top \mathbf{f}(x^i)$$

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

$\mathbf{w}^\top \mathbf{f}$

$\mathbf{w}_o^\top \mathbf{f}$

$\mathbf{w}_+^\top \mathbf{f}$

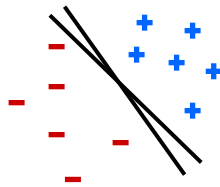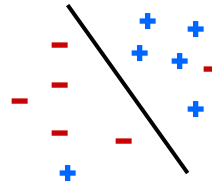$\mathbf{w}_-^\top \mathbf{f}$

# Separability

- A data set is (linearly) *separable* in a feature space if some linear classifier classifies all points correctly.
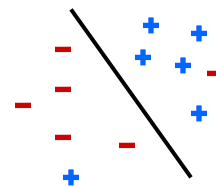
Separable               Non-Separable

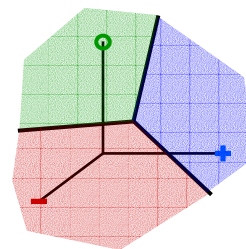- If a data set is separable, there are usually multiple separating hypotheses.

---

# Caution about Diagrams

- A diagram you'll often see:
    - Two-class classification
    - Fractional feature values
    - Mixed regions $\rightarrow$ non-separable
    - Sample complexity

- Common NLP case:
    - Multi-class classification
        - Each input corresponds to $|Y|$ points $f_i(y)$ (one per class)
    - (Mostly) 0/1 features
        - Data on the "corners"
    - Everything's separable
    - Coupon collection

# Linear Models: Naïve-Bayes

- (Multinomial) Naïve-Bayes: $\mathbf{x}^i = d_1, d_2, \cdots d_n$

$$\mathbf{f}_i(\mathbf{y}) = [\quad 0 \qquad 1, \qquad \#v_1, \qquad \#v_2, \qquad \cdots \qquad \#v_{|V|} \qquad 0 \quad]$$
$$\mathbf{w} = [\quad \cdots \quad \log P(y), \quad \log P(v_1|y), \quad \log P(v_2|y), \quad \cdots \quad \log P(v_n|y) \quad \cdots \quad]$$
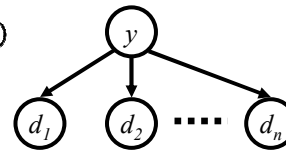
$$
\begin{aligned}
\text{score}(\mathbf{x}_i, \mathbf{y}, \mathbf{w}) &= \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) \\
&= \log P(\mathbf{y}) + \sum_k \#v_k \log P(v_k|\mathbf{y}) \\
&= \log \left( P(\mathbf{y}) \prod_k P(v_k|\mathbf{y})^{\#v_k} \right) \\
&= \log \left( P(\mathbf{y}) \prod_{d \in \mathbf{x}^i} P(d|\mathbf{y}) \right) \\
&= \log P(\mathbf{x}^i, \mathbf{y})
\end{aligned}
$$



---

# Bad Model Assumptions

**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:
- $P(r,+,+) = (\frac{1}{2})(\frac{3}{4})(\frac{3}{4})$
- $P(s,+,+) = (\frac{1}{2})(\frac{1}{4})(\frac{1}{4})$
- $P(r|+,+) = 9/10$
- $P(s|+,+) = 1/10$

# Worse Model Assumptions

### Reality

**Lights Working**



P(g,r,w) = 3/7        P(r,g,w) = 3/7

**Lights Broken**



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
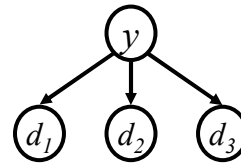
---

# Details: 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
- Discriminative models can partially compensate for wrong models
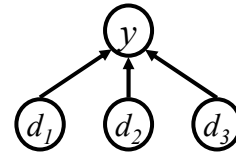
# Generative vs Discriminative

- Generative Models
  - Joint density over P(X,Y)
  - E.g. Naïve-Bayes, HMMs, PCFGs
  - Model assumptions allow decomposition into small factors which can be estimated independently
  - Do not set weights to account for feature interactions

- Discriminative Models
  - Predict Y given X, not always distributions
  - E.g. maximum entropy, SVMs, perceptrons
  - Set weights to account for feature interactions
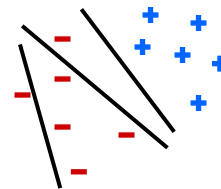  - Require *inference on training set* to evaluate hypotheses

# Linear Models: Perceptron

- Simple discriminative method for intuition

$$\mathbf{y}' = \arg\max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{\eta\left(\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y}')\right)}_{\Delta_i(\mathbf{y}')}$$

- This is a procedure, not an optimization problem!
  - May not converge if non-separable
  - Noisy

- Voted / averaged perceptron [Freund & Schapire 99, Collins 02]
  - Regularize / reduce variance by aggregating over iterations

# Objective Functions

- Reminder: $score(\mathbf{x}^i, \mathbf{y}, \mathbf{w}) = \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})$

- What do we want from weights?
  - Depends!
  - Minimize (training) errors?

$$\sum_i step\left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y} \neq \mathbf{y}^i} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})\right)$$

  - This is the "zero-one loss"
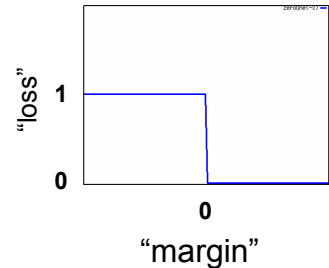    - Discontinuous, minimizing is NP-complete
    - Not really what we want anyway
  - Maxents and SVMs have losses related to the zero-one loss



"loss"

**1**

**0**

**0**

"margin"

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

---

# Linear Models: Maximum Entropy

- Maximum entropy (logistic regression)
  - Use the activations as probabilities:

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x},\mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x},\mathbf{y}'))}$$

⟵ Make positive

⟵ Normalize

  - Maximize the (log) conditional likelihood of training data

$$\max_{\mathbf{w}} \ \log \prod_i P(\mathbf{y}^i|\mathbf{x}^i, \mathbf{w}) = \sum_i \log\left(\frac{\exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i))}{\sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}))}\right)$$

$$\max_{\mathbf{w}} \ \sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}))\right)$$

"soft margin"

# "Soft-Max"

$$\max(a, b) \approx \log\left(\exp(a) + \exp(b)\right)$$

$$\max(a, b) \qquad \log\left(\exp(a) + \exp(b)\right)$$



# Maximum Entropy II

- Also: regularization (smoothing)

$$\max_{\mathbf{w}} \quad \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right) - k||w||^2$$

- Maximize likelihood = Minimize "log-loss"

$$\min_{\mathbf{w}} \quad k||w||^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

- Motivation
  - Connection to maximum entropy principle
  - Might want to do a good job of being uncertain on noisy cases…
  - … in practice, though, posteriors are pretty peaked

# Log-Loss

- If we view maxent as a minimization problem:

$$\min_{\mathbf{w}} \quad k\|w\|^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right)$$

- This minimizes the "log-loss" on each example
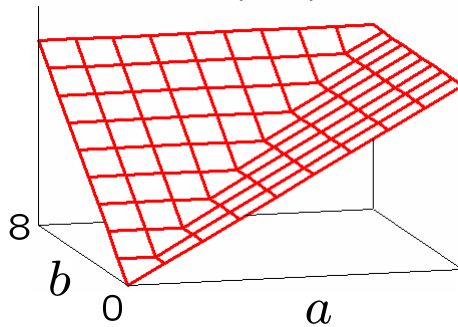
$$- \left[ \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})) \right]$$

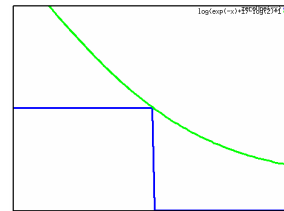$$- \log \left( \frac{\exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i))}{\sum_{\mathbf{y}} \exp(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}))} \right) = - \log P(\mathbf{y}^i | \mathbf{x}^i, \mathbf{w})$$

- Log-loss bounds zero-one loss

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

---

# SVMs

- SVM Try 1: Separate the training data

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

$$\mathbf{w}^\top \mathbf{f}(...\text{win election}..., POLITICS) \geq \mathbf{w}^\top \mathbf{f}(...\text{win election}..., SPORTS)$$

$$\mathbf{w}^\top \mathbf{f}(...\text{win election}..., POLITICS) \geq \mathbf{w}^\top \mathbf{f}(...\text{win election}..., OTHER)$$

1. This is an entire feasible space; need an objective function!

2. Training data may not even be separable

# Maximum Margin

- SVM Try 2: find the maximum margin separator

$$\max_{||\mathbf{w}||\leq 1} \gamma$$

$$\ell_i(\mathbf{y}) = \begin{cases} 0 & \text{if } \mathbf{y} = \mathbf{y}^i \\ 1 & \text{if } \mathbf{y} \neq \mathbf{y}^i \end{cases}$$

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \gamma \ell_i(\mathbf{y}) \quad \forall i, \forall \mathbf{y}$$

$\mathbf{w}^\top \mathbf{f}(\text{win election}, POLITICS) \geq \mathbf{w}^\top \mathbf{f}(\text{win election}, SPORTS) + \gamma$    1

$\mathbf{w}^\top \mathbf{f}(\text{win election}, POLITICS) \geq \mathbf{w}^\top \mathbf{f}(\text{win election}, OTHER) + \gamma$    1

$\mathbf{w}^\top \mathbf{f}(\text{win election}, POLITICS) \geq \mathbf{w}^\top \mathbf{f}(\text{win election}, POLITICS)$    0

$\ell$



---

# Why Max Margin?

- Why do this?  Various arguments:
  - Decisions on training points are maximally robust to "feature jitter"
  - As we'll see, solution depends only on the boundary cases, or *support vectors* (but remember how this diagram is broken!)
  - Sparse solutions (features not in support vectors get zero weight)
  - Generalization bound arguments



*Support vectors*

# Max Margin / Small Norm

- SVM Try 3: find the smallest w which separates data

Remember this condition? ➡️

$$\max_{||\mathbf{w}||\leq 1} \quad \gamma$$
$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \gamma \ell_i(\mathbf{y}) \quad \forall i, \forall \mathbf{y}$$

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

$$\min_{\mathbf{w}} \quad \frac{1}{2}||\mathbf{w}||^2$$
$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^*) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + 1\ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

---

# Max Gamma to Min W

$$\max_{||\mathbf{w}||\leq 1} \quad \gamma$$
$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \gamma \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$
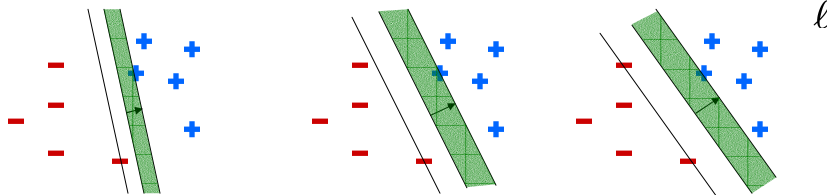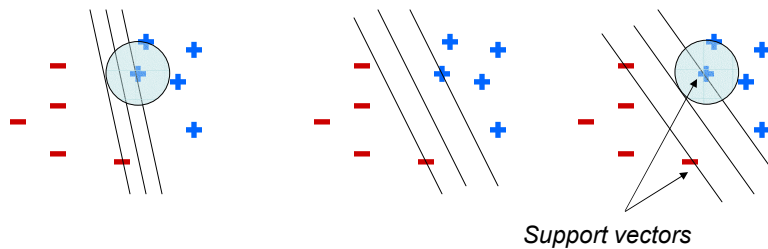
$$\mathbf{w} = \gamma u$$

$$\gamma = 1/||u||$$

$$\max_{||\gamma u||\leq 1} \quad 1/||u||^2$$
$$\text{s.t.} \quad \gamma u^\top \mathbf{f}_i(\mathbf{y}^i) \geq \gamma u^\top \mathbf{f}_i(\mathbf{y}) + \gamma \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

$$\max_{||\gamma u||\leq 1} \quad 1/||u||^2$$
$$\text{s.t.} \quad u^\top \mathbf{f}_i(\mathbf{y}^i) \geq u^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

$$\min_{||\gamma u||\geq 1} \quad ||u||^2$$
$$\text{s.t.} \quad u^\top \mathbf{f}_i(\mathbf{y}^i) \geq u^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

$$\min_{u} \quad ||u||^2$$
$$\text{s.t.} \quad u^\top \mathbf{f}_i(\mathbf{y}^i) \geq u^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

$$\min_{u} \quad \frac{1}{2}||u||^2$$
$$\text{s.t.} \quad u^\top \mathbf{f}_i(\mathbf{y}^i) \geq u^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

$$\min_{\mathbf{w}} \quad \frac{1}{2}||\mathbf{w}||^2$$
$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$
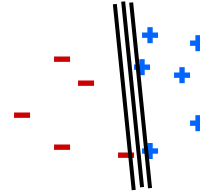
# Maximum Margin

- SVM Try 4: allow for non-separability
  - Add slack to the constraints
  - Make objective pay (linearly) for slack:

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

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

  - C is called the *capacity* of the SVM – the smoothing knob (more on this later)

- Learning:
  - Can stick this into Matlab if you want
  - Constrained optimization is hard; better methods!

---

# Min-Max Formulation

- We have a constrained minimization

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

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

- …but we can solve for $\xi_i$

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

$$\forall i, \quad \xi_i = \max_{\mathbf{y}} \left[ \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \right] - \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i)$$

- Giving

$$\min_{\mathbf{w}} \frac{1}{2}||\mathbf{w}||^2 - C \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}} \left[ \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \right] \right)$$

# Max vs "Soft-Max" Margin

- SVMs:

$$\min_{\mathbf{w}} k||\mathbf{w}||^2 - \sum_i \left( \underbrace{\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}} \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(y) \right)}_{\text{Hard (Penalized) Margin}} \right)$$

- Maxent:

$$\min_{\mathbf{w}} k||w||^2 - \sum_i \left( \underbrace{\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log \sum_{\mathbf{y}} \exp \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) \right)}_{\text{Soft Margin}} \right)$$

- 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"

---

# Hinge Loss

- Consider the per-instance SVM objective:

$$\min_{\mathbf{w}} k||\mathbf{w}||^2 - \sum_i \left( \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}} \left[ \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(y) \right] \right)$$

- This is called the "hinge loss"
  - Upper bounds zero-one 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



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

# Loss Functions: I

- **Zero-One Loss**

$$\sum_i step\left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}\neq\mathbf{y}^i} \mathbf{w}^\top \mathbf{f}_i(\mathbf{y})\right)$$

- **Hinge**

$$\sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}}\left[\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(y)\right]\right)$$

- **Log**

$$\sum_i \left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \log\sum_{\mathbf{y}} \exp\left(\mathbf{w}^\top \mathbf{f}_i(\mathbf{y})\right)\right)$$



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

# Loss Functions: II

# Loss Functions: III



# Outline

- Part I: Flat Classification
  - Linear classifiers and loss functions
  - Primal and dual SVM formulations
  - Training SVMs

- Part II: Structured Classification
  - Structured linear classifiers
  - Factored learning formulations
  - Experimental results

# Dual Formulation

- We want to optimize:

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

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

- This is hard because of the constraints.

- Solution: method of Lagrange multipliers

---

# 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 $\Lambda$ (we'll show this):

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

# Primal Game

- Original: $f(\mathbf{w}^*) \;=\; \min_{\mathbf{w}} \; f(\mathbf{w}) \quad s.t. \; g(\mathbf{w}) \geq 0$

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

  $\Lambda(\mathbf{w}) = \max_{\alpha \geq 0} \, [f(\mathbf{w}) - \alpha \, g(\mathbf{w})]$

  $f$

  $g > 0 \quad g = 0 \quad g < 0$

- Claim: primal game solves the original constrained problem:

  $$\min_{\mathbf{w}} \max_{\alpha \geq 0} \Lambda(\mathbf{w}, \alpha) = \min_{\mathbf{w}} \Lambda(\mathbf{w}) \;=\; f(\mathbf{w}^*)$$

- Proof: consider the value of

  $$\Lambda(\mathbf{w}) = \max_{\alpha \geq 0} \, [f(\mathbf{w}) - \alpha \, g(\mathbf{w})]$$

  $g(\mathbf{w}) = 0 \Rightarrow f(\mathbf{w})$

  $g(\mathbf{w}) > 0 \Rightarrow f(\mathbf{w})$

  $g(\mathbf{w}) < 0 \Rightarrow \infty$

  $\Lambda(\mathbf{w})$  $\boxed{f \quad \infty}$  $\Rightarrow$  $\min_{\mathbf{w}} \Lambda(\mathbf{w}) = \min_{\mathbf{w}:g \geq 0} f(\mathbf{w}) = f(\mathbf{w}^*)$


# Dual Game

- Original: $f(\mathbf{w}^*) \;=\; \min_{\mathbf{w}} \; f(\mathbf{w}) \quad s.t. \; g(\mathbf{w}) \geq 0$

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

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

- Claim: dual game also solves the original problem:

  $$\max_{\alpha \geq 0} \min_{\mathbf{w}} \Lambda(\mathbf{w}, \alpha) = \max_{\alpha \geq 0} \Lambda(\alpha) \;=\; f(\mathbf{w}^*)$$

- Proof:

  Case I: Constraint Inactive

  Case II: Constraint Active

# Dual Game IIa

- Lagrangian: $\Lambda(\alpha) = \min_{\mathbf{w}} \left[ f(\mathbf{w}) - \alpha\, g(\mathbf{w}) \right]$
- Claim: $\max_{\alpha \geq 0} \min_{\mathbf{w}} \Lambda(\mathbf{w}, \alpha) = \max_{\alpha \geq 0} \Lambda(\alpha) \;=\; f(\mathbf{w}^*)$

Case I: Constraint Inactive



At $\mathbf{w}^*$, $g > 0$, so if $\alpha > 0$,

$f(\mathbf{w}^*) - \alpha\, g(\mathbf{w}^*) < f(\mathbf{w}^*)$,

$\Lambda(\alpha) < f(\mathbf{w}^*)$

But $\Lambda(0) = f(\mathbf{w}^*)$

So $\max_{\alpha \geq 0} \Lambda(\alpha) = f(\mathbf{w}^*)$


# Dual Game IIb

- Lagrangian: $\Lambda(\alpha) = \min_{\mathbf{w}} \left[ f(\mathbf{w}) - \alpha\, g(\mathbf{w}) \right]$
- Claim: $\max_{\alpha \geq 0} \min_{\mathbf{w}} \Lambda(\mathbf{w}, \alpha) = \max_{\alpha \geq 0} \Lambda(\alpha) \;=\; f(\mathbf{w}^*)$

Case II: Constraint Active



At $\mathbf{w}^*$, $g = 0$, so $\forall \alpha$,

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

so $\forall \alpha$, $\Lambda(\alpha) < f(\mathbf{w}^*)$

At $\mathbf{w}^*$, $\nabla f \neq 0$, but

$\exists \alpha^*$ s.t. $\nabla f(\mathbf{w}^*) = \alpha^* \nabla g(\mathbf{w}^*)$

At $\alpha^*$, $\nabla \Lambda(\alpha^*, \mathbf{w}^*) = \nabla f - \alpha^* \nabla g = 0$

so $\Lambda(\alpha^*) = f(\mathbf{w}^*)$

# Lagrangian for SVMs

- Primal constrained problem:

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

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

- Lagrangian:

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

---

# Dual Formulation II

- Duality tells us that

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

has the same value as

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

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

# Dual Formulation III

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

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

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

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

$$\frac{\partial \Lambda(\mathbf{w}, \xi, \alpha)}{\partial \xi_i} = C - \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})$$

$$\frac{\partial \Lambda(\mathbf{w}, \xi, \alpha)}{\partial \xi_i} = 0 \implies \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y}) = C$$

# Dual Formulation IV

- We now know that for fixed $\alpha$, the minimum of

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

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

- Plugging these back into $\Lambda$:

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

# Dual Formulation V

- This doesn't reference the primal weights w at all, so we can now worry about the outer max problem:

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

$$\text{s.t.} \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C \quad \forall i$$

- And this solves the original constrained primal:

$$\max_{\alpha \geq 0} \Lambda(\alpha) = \max_{\alpha \geq 0} \min_{\mathbf{w},\xi} \Lambda(\mathbf{w},\xi,\alpha) \;=\; f(\mathbf{w}^*)$$

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

---

# What are the Alphas?

- Each example (and label) gave to a primal constraint

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

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}) \quad \forall i, \mathbf{y}$$

- 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 form the weights:

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

*Support vectors*

# Outline

- Part I: Flat Classification
  - Linear classifiers and loss functions
  - Primal and dual SVM formulations
  - Training SVMs

- Part II: Structured Classification
  - Structured linear classifiers
  - Factored learning formulations
  - Experimental results

# Back to Learning SVMs

- We want to find $\alpha$ which maximize

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

$$\text{s.t.} \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C \quad \forall i$$

- 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 Ascent I

- Consider the separable (soft-margin) SVM problem:

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

- In coordinate ascent, we maximize one variable at a time
- Despite all the mess, $Z$ is just a quadratic in each $\alpha_i(\mathrm{y})$



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

0     0

- If the unconstrained argmin on a coordinate is at a negative $\alpha$, just clip to zero!


# Coordinate Ascent II

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

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

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

- This approach works fine for the separable case

# Bi-Coordinate Descent I

- In the non-separable case, it's (a little) harder:

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

$$\text{s.t.} \quad \sum_{\mathbf{y}}\alpha_i(\mathbf{y}) = C \quad \forall i$$

- Here, we can't update just a single alpha, because of the sum-to-C constraints
- Instead, we can optimize two at once, shifting "mass" from one $\mathbf{y}$ to another:



---

# Bi-Coordinate Descent II

- Choose an example $i$, and two labels $\mathbf{y}_1$ and $\mathbf{y}_2$:

$$t = \frac{(\ell_i(\mathbf{y}_1) - \ell_i(\mathbf{y}_2)) - \left(\sum_{i,\mathbf{y}}\alpha_i(\mathbf{y})\left(\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y})\right)\right)^\top (\mathbf{f}_i(\mathbf{y}_2) - \mathbf{f}_i(\mathbf{y}_1))}{\|\mathbf{f}_i(\mathbf{y}_2) - \mathbf{f}_i(\mathbf{y}_1)\|^2}$$

$$\mathbf{y}_1 \to \min(\mathbf{y}_1 + t, \mathbf{y}_1 + \mathbf{y}_2)$$
$$\mathbf{y}_2 \to \max(\mathbf{y}_2 - t, 0)$$



$Z(t)$

- This is a sequential minimal optimization update, but it's not the same one as in [Platt 98]

# SMO

- Naïve SMO:

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

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

```
while (not converged) {
    visit each example i {
        for each pair of labels (y₁, y₂) {
            bi-coordinate-update(i, y₁, y₂)
        }
    }
}
```

- Time per iteration: $O(|x||\mathcal{Y}|^2)$
- Smarter SMO:
  - Can speed this up by being clever about skipping examples and label pairs which will make little or no difference

all examples

all label pairs

---

| DOCUMENT | FEATURES | FEATURE DELTAS | ALPHAS |
|---|---|---|---|
| win game | S-win, S-game | -- 0 -- | |
| | P-win, P-game | PW=1, SW=-1, PG=1,... | |
| | O-win, O-game | OW=1, SW=-1, OG=1,... | |
| win vote | S-win, S-vote | SW=1, PW=-1, SV=1,... | |
| | P-win, P-vote | -- 0 -- | |
| | O-win, O-vote | OW=1, PW=-1, OV=1,... | |
| movie | S-movie | SM=1, OM=-1 | |
| | P-movie | PM=1, OM=-1 | |
| | O-movie | -- 0 -- | |

WEIGHTS

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

-- 0 --

# Outline

- Part I: Flat Classification
  - Linear classifiers and loss functions
  - Primal and dual SVM formulations
  - Training SVMs

- Part II: Structured Classification
  - Structured linear classifiers
  - Factored learning formulations
  - Experimental results



"Don't worry, Howard. The big questions are multiple choice."

# Handwriting Recognition

x               y

 ➡️ **brace**

Sequential structure

---

# CFG Parsing

x               y

The screen was
a sea of red ➡️

```
                    S
              ┌─────┴─────┐
             NP           VP
          ┌───┴──┐    ┌────┴────┐
         DT  NN  VBD          NP
          │   │    │      ┌────┴────┐
        The screen was   NP        PP
                       ┌──┴──┐   ┌──┴──┐
                      DT  NN  IN      NP
                       │   │   │       │
                       a  sea  of     NN
                                       │
                                      red
```

Recursive structure

# Bilingual Word Alignment

**x**

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

**En vertu des nouvelles propositions, quel est le coût prévu de perception des droits?**

**y**

Combinatorial structure

---

# Structured Models

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

space of feasible outputs

Assumption:

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

Score = sum of local "part" scores

Parts = nodes, edges, productions

# Chain Markov Net (aka CRF*)

$$P(\mathbf{y} \mid \mathbf{x}) \propto \prod_j \underbrace{\phi(\mathbf{x}_j, y_j)} \prod_{jk} \underbrace{\phi(\mathbf{x}_{jk}, y_j, y_k)}$$

$$\phi(\mathbf{x}_j, y_j) = \exp\left\{\mathbf{w}_N^\top \mathbf{f}_N(\mathbf{x}_j, y_j)\right\} \qquad \text{N = Node}$$

$$\phi(\mathbf{x}_{jk}, y_j, y_k) = \exp\left\{\mathbf{w}_E^\top \mathbf{f}_E(\mathbf{x}_{jk}, y_j, y_k)\right\} \qquad \text{E = Edge}$$



$\mathbf{f}_E(\mathbf{x}_{jk}, y_j, y_k)$

$[\cdots I(y_j = \text{'z'}, y_k = \text{'a'}) \cdots]$

$\mathbf{f}_N(\mathbf{x}_j, y_j)$

$[\cdots I(\mathbf{x}_j[3,4] = 1, y_j = \text{'z'}) \cdots]$

*Lafferty et al. 01

---

# Chain Markov Net (aka CRF*)

$$P(\mathbf{y} \mid \mathbf{x}) \propto \prod_j \underbrace{\phi(\mathbf{x}_j, y_j)} \prod_{jk} \underbrace{\phi(\mathbf{x}_{jk}, y_j, y_k)} = \exp\left\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\right\}$$

$$\prod_j \phi(\mathbf{x}_j, y_j) = \exp\left\{\sum_j \mathbf{w}_N^\top \mathbf{f}_N(\mathbf{x}_j, y_j)\right\} = \exp\left\{\mathbf{w}_N^\top \mathbf{f}_N(\mathbf{x}, \mathbf{y})\right\}$$

$$\prod_{jk} \phi(\mathbf{x}_{jk}, y_j, y_k) = \exp\left\{\sum_{jk} \mathbf{w}_E^\top \mathbf{f}_E(\mathbf{x}_{jk}, y_j, y_k)\right\} = \exp\left\{\mathbf{w}_E^\top \mathbf{f}_E(\mathbf{x}, \mathbf{y})\right\}$$



$$\mathbf{f}_N(\mathbf{x}, \mathbf{y}) \equiv \sum_j \mathbf{f}_N(\mathbf{x}_j, y_j)$$

$$\mathbf{f}_E(\mathbf{x}, \mathbf{y}) \equiv \sum_{jk} \mathbf{f}_E(\mathbf{x}_{jk}, y_j, y_k)$$

$$\mathbf{f}(\mathbf{x}, \mathbf{y}) \equiv \begin{pmatrix} \mathbf{f}_N(\mathbf{x}, \mathbf{y}) \\ \mathbf{f}_E(\mathbf{x}, \mathbf{y}) \end{pmatrix} \quad \mathbf{w} \equiv \begin{pmatrix} \mathbf{w}_N \\ \mathbf{w}_E \end{pmatrix}$$

*Lafferty et al. 01

# CFG Parsing

$$P(\mathbf{y} \mid \mathbf{x}) \propto \prod_{A \to \alpha \in (\mathbf{x}, \mathbf{y})} \phi(A \to \alpha)$$

#(NP → DT NN)

...

$$\mathbf{f} : \mathcal{X} \times \mathcal{Y} \to \Re^d$$

#(PP → IN NP)

...

#(NN → 'sea')

$$\prod_{A \to \alpha \in (\mathbf{x}, \mathbf{y})} \exp\left\{\mathbf{w}^\top \mathbf{f}(A \to \alpha)\right\} = \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\}$$

---

# Bilingual Word Alignment

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

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

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

k

$y_{jk}$

j

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

- association
- position
- orthography

# Probabilistic Alignment?

$$P(\mathbf{y} \mid \mathbf{x}) = \frac{\exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})\}}{\sum_{\mathbf{y}'} \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}')\}}$$

← *#P-Complete*

Need to sum over all possible matchings

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

$y_{jk}$

k

j

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

---

# OCR Example

- We want:

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

- Equivalently:

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) \ > \ \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaaa"})$$

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) \ > \ \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"aaaab"})$$

...

$$\mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"brace"}) \ > \ \mathbf{w}^\top \mathbf{f}(\text{brace}, \text{"zzzzz"})$$

**a lot!**

# Parsing Example

- We want:

$$\arg\max_{\mathbf{y}} \ \mathbf{w}^{\top}\mathbf{f}(\ \text{'It was red'}\ ,\mathbf{y}\ ) \ = \ \begin{smallmatrix}S\\A\ B\\C\ D\end{smallmatrix}$$

- Equivalently:

$$\mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\A\ B\\C\ D\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\A\ B\\D\ F\end{smallmatrix})$$

$$\mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\A\ B\\C\ D\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\A\ B\\C\ D\end{smallmatrix})$$

...

$$\mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\A\ B\\C\ D\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\text{'It was red'},\ \begin{smallmatrix}S\\E\ F\\G\ H\end{smallmatrix})$$

**a lot!**

---

# Alignment Example

- We want:

$$\arg\max_{\mathbf{y}} \ \mathbf{w}^{\top}\mathbf{f}(\ \begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix}\ ,\mathbf{y}\ ) \ = \ \begin{smallmatrix}1 - 1\\2 - 2\\3 - 3\end{smallmatrix}$$

- Equivalently:

$$\mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2-2\\3-3\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2\times 2\\3-3\end{smallmatrix})$$

$$\mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2-2\\3-3\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2\times 2\\3-3\end{smallmatrix})$$

...

$$\mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2-2\\3-3\end{smallmatrix}) \ > \ \mathbf{w}^{\top}\mathbf{f}(\begin{smallmatrix}\text{'What is the'}\\\text{'Quel est le'}\end{smallmatrix},\ \begin{smallmatrix}1-1\\2\times 2\\3-3\end{smallmatrix})$$

**a lot!**

# Structured Loss

| b | ✗ | a | ✗ | e | 2 |
|---|---|---|---|---|---|
| b | r | ✗ | ✗ | e | 2 |
| b | r | ✗ | c | e | 1 |
| b | r | a | c | e | 0 |

| 0 | 1 | 2 | 3 | | 0 | 1 | 2 | 2 |

'It was red'    'What is the'  'Quel est le'

---

# Max Margin Estimation

- Given training example $\mathbf{x}^i, \mathbf{y}^i$ we want:

$$\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) > \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) \quad \forall i, \mathbf{y} \neq \mathbf{y}^i$$

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

- Maximize loss weighted margin:  $\gamma \ell_i(\mathbf{y})$

$$\ell_i(\mathbf{y}) = \sum_j I(y_j^i \neq y_j) \qquad \text{\# of mistakes in } \mathbf{y}$$

*Collins 02, Altun et al 03, Taskar 03

# Large margin estimation

- **Brute force enumeration**

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

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}), \quad \forall i, \mathbf{y}$$

- **Min-max formulation**

$$\min_{\mathbf{w}} \quad \frac{1}{2}||\mathbf{w}||^2 - C\left(\sum_i \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}}\left[\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})\right]\right)$$

- Plug-in linear program for loss-augmented inference

$$\max_{\mathbf{y}}\left[\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})\right]$$

---

# Min-max formulation

$$\max_{\mathbf{y}}\left[\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})\right]$$

Assume linear loss (Hamming):
$$\ell_i(\mathbf{y}) = \sum_p \ell_{i,p}(\mathbf{y}_p)$$

DP Inference
$$\max_{\mathbf{y}}\left[\sum_p \mathbf{w}^\top \mathbf{f}(\mathbf{x}_p, \mathbf{y}_p) + \ell_{i,p}(\mathbf{y}_p)\right]$$

LP inference
$$\max_{\substack{\mathbf{z}\geq 0; \\ \mathbf{A}\mathbf{z}\leq\mathbf{b};}} \mathbf{q}^\top \mathbf{z}$$

# $\mathbf{y} \Rightarrow \mathbf{z}$ Map for Markov Nets

$$\mathbf{y} = \text{'ababb'}$$

|   | $z_1(m)$ | $z_2(m)$ | $z_3(m)$ | $z_4(m)$ | $z_5(m)$ |
|---|---|---|---|---|---|
| a | 1 | 0 | 1 | 0 | 0 |
| b | 0 | 1 | 0 | 1 | 1 |
| : | : | : | : | : | : |
| z | 0 | 0 | 0 | 0 | 0 |

$z_{12}(m,n)$  $z_{23}(m,n)$  $z_{34}(m,n)$  $z_{45}(m,n)$

|   | a | b | . | z | a | b | . | z | a | b | . | z | a | b | . | z |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | 0 | 1 | . | 0 | 0 | 0 | . | 0 | 0 | 1 | . | 0 | 0 | 0 | . | 0 |
| b | 0 | 0 | . | 0 | 1 | 0 | . | 0 | 0 | 0 | . | 0 | 0 | 1 | . | 0 |
| : | . | . | . | 0 | . | . | . | 0 | . | . | . | 0 | . | . | . | 0 |
| z | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| a | b | . | z | a | b | . | z | a | b | . | z | a | b | . | z |

---

# Markov Net Inference LP

$$\max_{\mathbf{z}} \quad \sum_{j,m} z_j(m) \left[ \mathbf{w}^\top \mathbf{f}_{\mathsf{N}}(\mathbf{x}_j, m) + \ell_j(m) \right]$$

$$+ \sum_{jk,m,n} z_{jk}(m,n) \left[ \mathbf{w}^\top \mathbf{f}_{\mathsf{E}}(\mathbf{x}_{jk}, m, n) + \ell_{jk}(m,n) \right] \left.\right\} \begin{array}{l} \mathbf{q}^\top \mathbf{z} \\ \mathbf{q} = \mathbf{F}^\top \mathbf{w} + \ell \end{array}$$

$z_k(n)$

| 0 | 1 | 0 | 0 |
|---|---|---|---|

$z_j(m)$

$$z_j(m) \geq 0; \qquad z_{jk}(m,n) \geq 0;$$

| 0 | | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|
| 0 | | 0 | 0 | 0 | 0 |
| 1 | | 0 | 1 | 0 | 0 |
| 0 | | 0 | 0 | 0 | 0 |

normalization $\quad \sum_m z_j(m) = 1$

agreement $\quad \sum_n z_{jk}(m,n) = z_j(m)$

$\left.\right\} \; A\mathbf{z} = \mathbf{b}$

$z_{jk}(m,n)$

Has integral solutions **z** for chains, trees

# CFG Chart



$z_{027}(\text{S}, \text{NP}, \text{VP})$

$z_{35}(\text{NP})$

- CNF tree = set of two types of parts:
  - Constituents $(A, s, e)$
  - CF-rules $(A \rightarrow B\,C, s, m, e)$

$$\mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{p \in \mathbf{y}} \mathbf{f}(\mathbf{x}, p)$$

---

# CFG Inference LP

$$\max_{\mathbf{z}} \quad \sum_{\substack{s<m<e \\ A \to B\,C}} z_{sme}(ABC) \left[ \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{sme}, ABC) + \ell_{sme}(ABC) \right] \left.\begin{array}{l} \mathbf{q}^\top \mathbf{z} \\ \mathbf{q} = \mathbf{F}^\top \mathbf{w} + \ell \end{array}\right.$$

$$\text{s.t.} \quad z_{se}(A) \geq 0 \qquad z_{sme}(ABC) \geq 0$$

root $\quad \displaystyle\sum_A z_{0,n}(A) = 1$

inside $\quad z_{se}(A) = \displaystyle\sum_{\substack{s<m<e \\ B,C}} z_{sme}(A, B, C)$

outside $\quad z_{se}(A) = \displaystyle\sum_{\substack{e<m\leq n \\ B,C}} z_{sme}(B, A, C) + \sum_{\substack{0\leq m<s \\ B,C}} z_{sme}(B, C, A)$

$\left.\rule{0pt}{3.5em}\right\} A\mathbf{z} = \mathbf{b}$

Has integral solutions **z**

# Matching Inference LP

$$\max_{\mathbf{z}} \quad \sum_{jk} z_{jk} \left[ \mathbf{w}^\top \mathbf{f}(\mathbf{x}_{jk}) + \ell_{jk} \right] \left.\right\} \begin{array}{l} \mathbf{q}^\top \mathbf{z} \\ \mathbf{q} = \mathbf{F}^\top \mathbf{w} + \ell \end{array}$$

$$\text{s.t.} \quad z_{jk} \geq 0$$

$$\text{degree} \quad \left. \begin{array}{l} \sum_k z_{jk} \leq 1 \\ \sum_j z_{jk} \leq 1 \end{array} \right\} A\mathbf{z} \leq \mathbf{b}$$

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

k

$z_{jk}$

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

Has integral solutions **z**

---

# LP Duality Recap

- Linear programming duality
  - Variables $\Rightarrow$ constraints
  - Constraints $\Rightarrow$ variables
- Optimal values are the same
  - When both feasible regions are bounded

$$\begin{array}{ll} \max_{\mathbf{z}} & \mathbf{c}^\top \mathbf{z} \\ \text{s.t.} & \mathbf{A}\mathbf{z} \leq \mathbf{b}; \\ & \mathbf{z} \geq 0. \end{array} \quad \Longleftrightarrow \quad \begin{array}{ll} \min_{\lambda} & \mathbf{b}^\top \lambda \\ \text{s.t.} & \mathbf{A}^\top \lambda \geq \mathbf{c}; \\ & \lambda \geq 0. \end{array}$$

# Min-max formulation

$$\min_{\mathbf{w}} \quad \frac{1}{2}\|\mathbf{w}\|^2 - C\left(\sum_i \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \max_{\mathbf{y}}\left[\mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y})\right]\right)$$

$$\max_{\substack{\mathbf{A}_i \mathbf{z}_i \leq \mathbf{b}_i \\ \mathbf{z}_i \geq 0}} \mathbf{q}_i^\top \mathbf{z}_i \qquad \Longleftrightarrow \qquad \min_{\substack{\mathbf{A}_i^\top \lambda_i \geq \mathbf{q}_i \\ \lambda_i \geq 0}} \mathbf{b}_i^\top \lambda_i$$

LP duality

$$\min_{\mathbf{w},\lambda} \quad \frac{1}{2}\|\mathbf{w}\|^2 - C\left(\sum_i \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \mathbf{b}_i^\top \lambda_i\right)$$

$$\text{s.t.} \quad \mathbf{A}_i^\top \lambda_i \geq \mathbf{q}_i; \quad \lambda_i \geq 0$$

$$\mathbf{q}_i = \mathbf{F}_i^\top \mathbf{w} + \ell_i$$

# Min-max formulation summary

$$\min_{\mathbf{w},\lambda} \quad \frac{1}{2}\|\mathbf{w}\|^2 - C\left(\sum_i \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \mathbf{b}_i^\top \lambda_i\right)$$

$$\text{s.t.} \quad \mathbf{A}_i^\top \lambda_i \geq \mathbf{F}_i^\top \mathbf{w} + \ell_i; \quad \lambda_i \geq 0, \ \forall i.$$

- Formulation produces concise QP for
  - Low-treewidth Markov networks
  - Context free grammars
  - Bipartite matchings
  - Many other problems with compact LP inference

*Taskar et al 04

# Factored Primal/Dual

$$\min_{\mathbf{w}, \lambda} \quad \frac{1}{2}||\mathbf{w}||^2 - C\left(\sum_i \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) - \mathbf{b}_i^\top \lambda_i\right)$$

$$\text{s.t.} \quad \mathbf{A}_i^\top \lambda_i \geq \mathbf{F}_i^\top \mathbf{w} + \ell_i; \quad \lambda_i \geq 0, \quad \forall i.$$

By QP duality $\quad \mathbf{w} = \sum_i C\mathbf{f}_i(\mathbf{y}^i) - \mathbf{F}_i \mu_i$

$$\max_{\mu} \quad \sum_i \ell_i^\top \mu_i - \frac{1}{2}\left\|\sum_i C\mathbf{f}_i(\mathbf{y}^i) - \mathbf{F}_i \mu_i]\right\|^2$$

$$\text{s.t.} \quad \mathbf{A}_i \mu_i \leq C\mathbf{b}_i; \quad \mu_i \geq 0, \quad \forall i.$$

Dual inherits structure from problem-specific inference LP

Variables $\boldsymbol{\mu}$ correspond to a decomposition of $\boldsymbol{\alpha}$ variables of the flat case

# Unfactored Primal/Dual

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

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}^i) + \xi_i \geq \mathbf{w}^\top \mathbf{f}_i(\mathbf{y}) + \ell_i(\mathbf{y}), \quad \forall i, \mathbf{y}$$

By QP duality $\quad \mathbf{w} = \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})[\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y})]$

$$\max_{\alpha} \quad \sum_{i,\mathbf{y}} \ell_i(\mathbf{y})\alpha_i(\mathbf{y}) - \frac{1}{2}\left\|\sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})[\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y})]\right\|^2$$

$$\text{s.t.} \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \quad \alpha_i \geq 0, \quad \forall i.$$
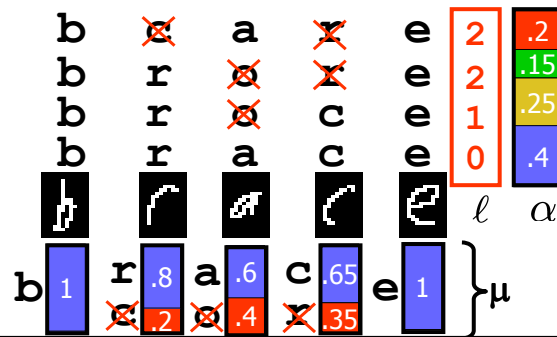
Exponentially many constraints/variables

# The Connection

$$\max_{\mu} \quad \sum_i \ell_i^\top \mu_i - \frac{1}{2} \left\| \sum_i C \mathbf{f}_i(\mathbf{y}^i) - \mathbf{F}_i \mu_i \right\|^2$$

$$\text{s.t.} \quad \mathbf{A}_i \mu_i = C \mathbf{b}_i; \quad \mu_i \geq 0, \quad \forall i.$$

$$\max_{\alpha} \quad \sum_{i,\mathbf{y}} \ell_i(\mathbf{y}) \alpha_i(\mathbf{y}) - \frac{1}{2} \left\| \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})[\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y})] \right\|^2$$

$$\text{s.t.} \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \quad \alpha_i \geq 0, \quad \forall i.$$

| b | ~~c~~ | a | ~~r~~ | e | 2 | .2 |
| b | r | ~~o~~ | ~~x~~ | e | 2 | .15 |
| b | r | ~~x~~ | c | e | 1 | .25 |
| b | r | a | c | e | 0 | .4 |

ℓ   α

**b** `1`   **r** `.8` **a** `.6` **c** `.65` **e** `1`  } µ
~~c~~ `.2`  ~~o~~ `.4`  ~~x~~ `.35`

---

# Structured SMO

$$\max_{\mu} \quad \sum_i \ell_i^\top \mu_i - \frac{1}{2} \left\| \sum_i C \mathbf{f}_i(\mathbf{y}^i) - \mathbf{F}_i \mu_i \right\|^2$$

$$\text{s.t.} \quad \mathbf{A}_i \mu_i = C \mathbf{b}_i; \quad \mu_i \geq 0, \quad \forall i.$$

$$\max_{\alpha} \quad \sum_{i,\mathbf{y}} \ell_i(\mathbf{y}) \alpha_i(\mathbf{y}) - \frac{1}{2} \left\| \sum_{i,\mathbf{y}} \alpha_i(\mathbf{y})[\mathbf{f}_i(\mathbf{y}^i) - \mathbf{f}_i(\mathbf{y})] \right\|^2$$

$$\text{s.t.} \quad \sum_{\mathbf{y}} \alpha_i(\mathbf{y}) = C; \quad \alpha_i \geq 0, \quad \forall i.$$

$\mu$ → **select & lift** → $\mathbf{y}'$ $\alpha$ $\mathbf{y}''$ → **SMO update** → $\mathbf{y}'$ $\alpha'$ $\mathbf{y}''$ → **project** → $\mu'$

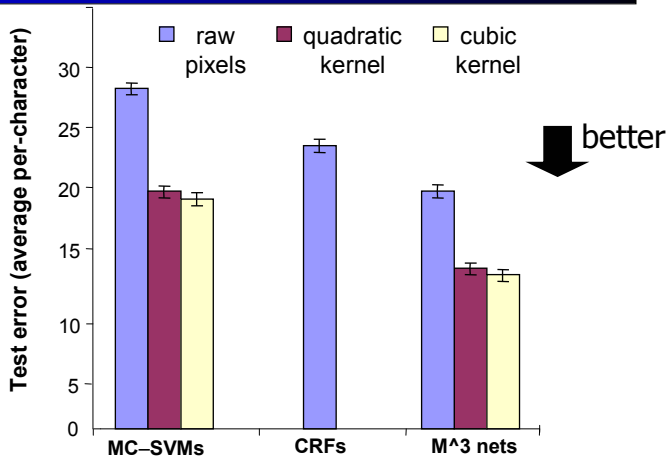*Taskar 04

# Outline

- Part I: Flat Classification
  - Linear classifiers and loss functions
  - Primal and dual SVM formulations
  - Training SVMs

- Part II: Structured Classification
  - Structured linear classifiers
  - Factored learning formulations
  - Experimental results
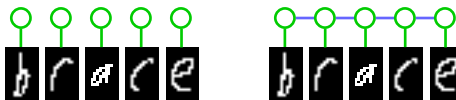
---

# Handwriting Recognition

Length: ~8 chars
Letter: 16x8 pixels
10-fold Train/Test
5000/50000 letters
600/6000 words

Models:
  Multiclass-SVMs
  CRFs
  $M^3$ nets

*Taskar et al 03

# Experimental Setup

- Standard Penn treebank split  (2-21/22/23)
- Generative baselines
    - Klein & Manning 03 and Collins 99
- Discriminative
    - Basic = max-margin version of K&M 03
    - Lexical & Lexical + Aux
- Lexical features (on constituent parts only)

$$t_{s-1} \; [t_s \; \ldots \; t_e] \; t_{e+1}$$
$$x_{s-1} \; [x_s \; \ldots \; x_e] \; x_{e+1}$$

$\leftarrow$  predicted tags

- Auxillary features
    - Flat classifier using same features
    - Prediction of K&M 03 on each span

---

# Results for sentences ≤40 words

| Model | LP | LR | $F_1$ |
|---|---|---|---|
| Generative | 86.37 | 85.27 | 85.82 |
| Lexical+Aux* | **87.56** | **86.85** | **87.20** |
| Collins 99* | 85.33 | 85.94 | 85.73 |

*Trained only on sentences ≤20 words

*Taskar et al 04

# Example

*The Egyptian president said he would visit Libya today to resume the talks.*
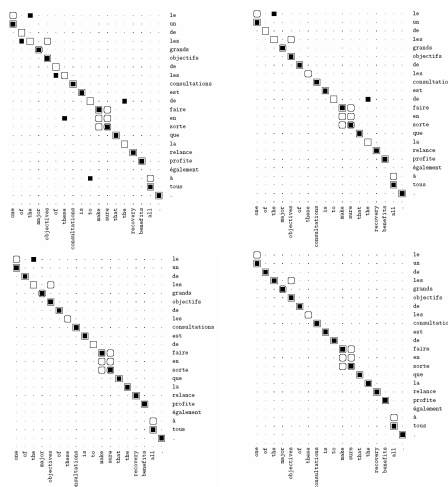
Generative model: Libya today is base NP

Lexical model: today is a one word constituent

---

# Word Alignment Results

- Hansards, 2M unlabeled, 100 labeled sentences

| Model | AER |
|---|---|
| Dice | 36.0 |
| IBM 4 | 9.7 |
| MM-Dice | 29.8 |
| +Distance | 17.2 |
| +Shape/Freq | 14.3 |
| +Next/Common | 9.6 |

## Generative/Discriminative Trade-offs

- Inference on training:
  - Discriminative methods require (repeated) inference on the training set, over the domains where the parameters interact
  - Generative models are primarily estimated from statistics of the training set (counting)
  - Inference can be much, much slower than counting

- Accounting for interactions:
  - Discriminative estimates take into account feature interactions, non-independence (note that conjunctive features are required to actually model interactions)

- Bias / variance
  - Discriminative methods tend to have higher variance, generative ones tend to have higher bias – but in general the discriminative techniques win on accuracy if properly regularized

## Likelihood/Margin Trade-offs

- Same as maxent vs. SVMs:
  - Sparse solutions, robust to "feature jitter"
  - Margin-based training often more accurate when posteriors are not needed

- Plus: unnormalizable models
  - For some models (e.g., matchings and a subclass of Markov networks), margin is tractable, likelihood is not!

# Conclusions

- **Today's tutorial:**
  - Flat SVMs from scratch
    - Objective functions and properties
    - Primal and dual formulations
    - How to learn them
  - Structured max-margin models
    - Concise, factored form
    - Efficient algorithms, strong empirical results
    - Applications: sequences, trees, matchings
- **Coming soon:**
  - Sequence modeling toolkit including M3Ns

  http://www.cs.berkeley.edu/~klein
  http://www.cs.berkeley.edu/~taskar

# References

Y. Altun, I. Tsochantaridis, and T. Hofmann. *Hidden Markov support vector machines.* ICML03.

M. Collins. *Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms*. EMNLP02

K. Crammer and Y. Singer. *On the algorithmic implementation of multiclass kernel-based vector machines.* JMLR01

J. Lafferty, A. McCallum, and F. Pereira. *Conditional random fields: Probabilistic models for segmenting and labeling sequence data.* ICML04

B. Taskar, C. Guestrin, D. Koller. *Maximum margin Markov Networks*. NIPS*03*

B. Taskar, D. Klein, M. Collins, D. Koller, C. Manning. *Maximum margin Parsing*. EMNLP*04*

B. Taskar. *Learning structured prediction models: a large margin approach. Stanford Univ. Thesis, 2004*