## SAILING LAB
Laboratory for Statistical Artificial InteLligence & INtegrative Genomics

# On Learning Sparse Structured Input-Output Models

# Eric Xing

epxing@cs.cmu.edu

Machine Learning Dept./Language Technology Inst./Computer Science Dept.
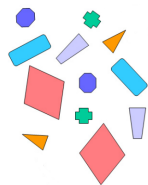
## Carnegie Mellon University

EMNLP 2012　　　　　　　　　　　　　　　　　　　　　　　　7/17/12　　**1**

---

## SAILING LAB
Laboratory for Statistical Artificial InteLligence & INtegrative Genomics

# Unstructured Prediction Problem



$$\mathbf{x} = \left( \begin{array}{ccc} x_{11} & x_{12} & \dots \end{array} \right) \quad \Rightarrow \quad \mathbf{y} = y_1$$

EMNLP 2012　　　　　　　　　　　　　　　　　　　　　　　　**2**

# Classical Predictive Models

- Input and output space: $\mathcal{X} \triangleq \mathbb{R}^{M_x} \qquad \mathcal{Y} \triangleq \{-1, +1\}$

- Predictive function $h(\mathbf{x})$: $\quad y^\star = h(\mathbf{x}) \triangleq \arg\max_{y \in \mathcal{Y}} F(\mathbf{x}, y; \mathbf{w})$

- Examples: $\qquad\qquad F(\mathbf{x}, y; \mathbf{w}) = g(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y))$

- Learning: $\qquad\qquad \hat{\mathbf{w}} = \arg\min_{\mathbf{w} \in \mathcal{W}} \ell(\mathbf{x}, y; \mathbf{w}) + \lambda R(\mathbf{w})$

  where $\ell(\cdot)$ represents a **convex loss**, and $R(\mathbf{w})$ is a **regularizer** preventing overfitting

| – Logistic Regression | – Support Vector Machines (SVM) |
|---|---|
| • Max-likelihood (or MAP) estimation | • Max-margin learning |
| $\max_{\mathbf{w}} \mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \sum_{i=1}^{N} \log p(y^i \mid \mathbf{x}^i; \mathbf{w}) + \mathcal{N}(\mathbf{w})$ | $\min_{\mathbf{w}, \xi} \frac{1}{2}\mathbf{w}^\top\mathbf{w} + C\sum_{i=1}^{N} \xi_i;$ |
| | s.t. $\forall i, \forall y' \neq y^i : \mathbf{w}^\top \Delta \mathbf{f}_i(y') \geq 1 - \xi_i, \ \xi_i \geq 0.$ |
| $\ell_{LL}(\mathbf{x}, y; \mathbf{w}) \triangleq \ln \sum_{y' \in \mathcal{Y}} \exp\{\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y')\} - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y)$ | $\ell_{MM}(\mathbf{x}, y; \mathbf{w}) \triangleq \max_{y' \in \mathcal{Y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y') - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, y) + \ell'(y', y)$ |

EMNLP 2012

3

---

# From Unstructured to Structured Prediction

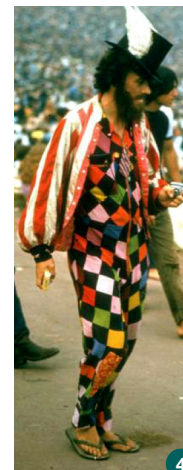- **Binary** classification: black-and-white decisions

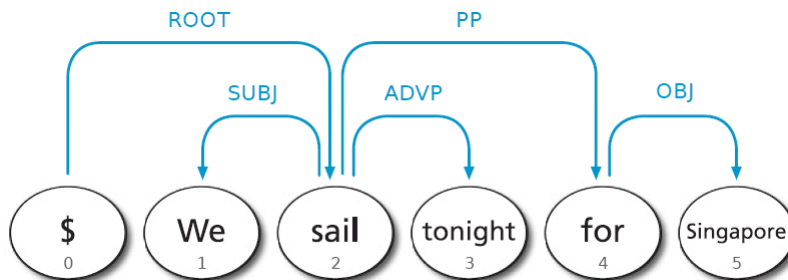- **Multi-class** classification: the world of technicolor

  - can be reduced to several binary decisions, but...
  - often better to handle multiple classes directly
  - how many classes? 2? 5? exponentially many?

- **Structured** prediction: many classes, strongly interdependent

  - Example: image segmentation (number of classes exponential to the # of segments)

$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix} \implies \mathbf{y} = \begin{pmatrix} y_{11} & y_{12} & \cdots \\ y_{21} & y_{22} & \cdots \\ \vdots & \vdots & \cdots \end{pmatrix}$$

EMNLP 2012

4

# Example I: Dependency Parsing of Sentences

ROOT                                    PP

SUBJ          ADVP                OBJ

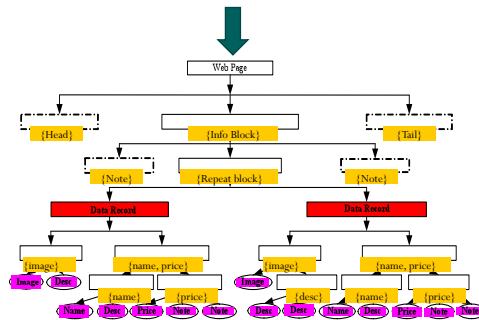$ 0    We 1    sail 2    tonight 3    for 4    Singapore 5

Challenge:
Structured outputs, and globally constrained to be a valid tree

# Example II: Text Summarization

*Australian novelist Peter Carey was awarded the coveted Booker Prize for fiction Tuesday night for his love story, "Oscar and Lucinda".*
*A panel of five judges unanimously announced the award of the $26,250 prize after an 80-minute deliberation during a banquet at London's ancient Guildhall.*
*The judges made their selection from 102 books published in Britain in the past 12 months and which they read in their homes.*
*Carey, who lives in Sydney with his wife and son, said in a brief speech that like the other five finalists he had been asked to attend with a short speech in his pocket in case he won.*

# Example III: Web-Data Extraction

7

# Example IV: Topic Discovery/Extraction

| "Arts" | "Budgets" | "Children" | "Education" |
|--------|-----------|------------|-------------|
| NEW | MILLION | CHILDREN | SCHOOL |
| FILM | TAX | WOMEN | STUDENTS |
| SHOW | PROGRAM | PEOPLE | SCHOOLS |
| MUSIC | BUDGET | CHILD | EDUCATION |
| MOVIE | BILLION | YEARS | TEACHERS |
| PLAY | FEDERAL | FAMILIES | HIGH |
| MUSICAL | YEAR | WORK | PUBLIC |
| BEST | SPENDING | PARENTS | TEACHER |
| ACTOR | NEW | SAYS | BENNETT |
| FIRST | STATE | FAMILY | MANIGAT |
| YORK | PLAN | WELFARE | NAMPHY |
| OPERA | MONEY | MEN | STATE |
| THEATER | PROGRAMS | PERCENT | PRESIDENT |
| ACTRESS | GOVERNMENT | CARE | ELEMENTARY |
| LOVE | CONGRESS | LIFE | HAITI |

The William Randolph Hearst Foundation will give $1.25 million to Lincoln Center, Metropolitan Opera Co., New York Philharmonic and Juilliard School. "Our board felt that we had a real opportunity to make a mark on the future of the performing arts with these grants an act every bit as important as our traditional areas of support in health, medical research, education and the social services," Hearst Foundation President Randolph A. Hearst said Monday in announcing the grants. Lincoln Center's share will be $200,000 for its new building, which will house young artists and provide new public facilities. The Metropolitan Opera Co. and New York Philharmonic will receive $400,000 each. The Juilliard School, where music and the performing arts are taught, will get $250,000. The Hearst Foundation, a leading supporter of the Lincoln Center Consolidated Corporate Fund, will make its usual annual $100,000 donation, too.

8

4

## Slide 9

# Structured Prediction Graphical Models

- Input and output space: $\mathcal{X} \triangleq \mathbb{R}_{X_1} \times, \ldots, \mathbb{R}_{X_K} \quad \mathcal{Y} \triangleq \mathbb{R}_{Y_1} \times, \ldots, \mathbb{R}_{Y_{K'}}$
- Convex loss function

- Conditional Random Fields (CRFs) (Lafferty et al 2001)
  - Based on a **Logistic Loss** (LR)
  - Max-likelihood estimation (point-estimate)

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \log \sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}')) \\ - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$
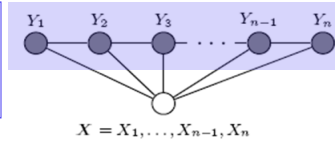
- Max-margin Markov Networks ($M^3Ns$) (Taskar et al 2003)
  - Based on a **Hinge Loss** (SVM)
  - Max-margin learning (point-estimate)

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \log \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}') \\ - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \ell(\mathbf{y}', \mathbf{y})$$

- Markov properties are encoded in the **feature functions** $\mathbf{f}(\mathbf{x}, \mathbf{y})$
  $$F(\mathbf{x}, y; \mathbf{w}) = g(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y))$$

$X = X_1, \ldots, X_{n-1}, X_n$

EMNLP 2012

**9**

---

## Slide 10

# Structured Prediction Graphical Models

- Conditional Random Fields (CRFs) (Lafferty et al 2001)
  - Based on a **Logistic Loss** (LR)
  - Max-likelihood estimation (point-estimate)

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \log \sum_{\mathbf{y}'} \exp(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}')) \\ - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$$
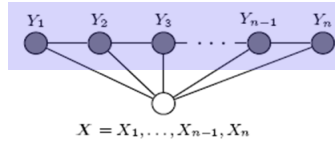
- Max-margin Markov Networks ($M^3Ns$) (Taskar et al 2003)
  - Based on a **Hinge Loss** (SVM)
  - Max-margin learning (point-estimate)

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) \triangleq \log \max_{\mathbf{y}'} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}') \\ - \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) + \ell(\mathbf{y}', \mathbf{y})$$

**Challenges:**
- **SPARSE "Interpretable"** prediction model
- **Prior** information of structures
- **Latent** structures/variables
- **Time** series and non-stationarity
- **Scalable** to large-scale problems (e.g., $10^4$ or larger input/output dimension)

$X = X_1, \ldots, X_{n-1}, X_n$

EMNLP 2012

**10**

# Main Claims

- The **sparse structures** of natural language data (input) and of the NLP tasks (output) can be utilized to improve the **quality** of the solution and **interpretability** of the solution

- Over-parameterized models such as conventional NB/LR/ SVM-style classifiers or parsers, topic models, or the related spectrum methods are not benefiting from sparse structures

- It is desirable to explore model spaces with structured sparsity for both **predictive models** (e.g., classifiers, parsers) and **explorative models** (e.g., topic models)

# Outline

- Sparse Structured Input-Output Models
  - … supervised learning
  - … convex optimization and log loss
  - … Frequentist-style shrinkage via regularization

- Sparse Topic Models
  - … unsupervised learning
  - … non-convex and likelihood-driven
  - … Bayesian-style posterior inference

- Sparse and Discriminative Topic Models?
  - … toward jointly explorative and predictive learning

## Slide 13

SAILING LAB

# Basic text classification

Class Label     Word counts     Feature strength

| 0 | learning |
|---|---|
| 3 | journal |
| 2 | intelligence |
| 0 | text |
| 0 | agent |
| 1 | internet |
| 0 | webwatcher |
| 0 | perlS |
| . | |
| . | |
| . | |
| 1 | volume |

1     =     [word counts]    X    **?**

$$y \quad = \quad X \quad \times \quad \beta$$

---

## Slide 14

SAILING LAB

# Basic text classification

Class Label     Word counts     Feature strength

| 0 | learning |
|---|---|
| 3 | journal |
| 2 | intelligence |
| 0 | text |
| 0 | agent |
| 1 | internet |
| 0 | webwatcher |
| 0 | perlS |
| . | |
| . | |
| . | |
| 1 | volume |

1     =     [word counts]    X

$$\beta^* = \arg\min_{\beta}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta)$$

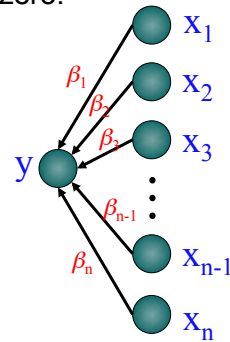Many non-zero coefficients:
Which words are truly significant?

# Sparsity: In a mathematical sense

- Consider least squares linear regression problem:
- Sparsity means most of the beta's are zero.

$$\hat{\boldsymbol{\beta}} = \mathrm{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

subject to:

$$\sum_{j=1}^{p} \mathbb{I}[|\beta_j| > 0] \leq C$$



- But this is not convex!!! Many local optima, computationally intractable.

---

# L1 Regularization (LASSO) (Tibshirani, 1996)

- A convex relaxation.

**Constrained Form**

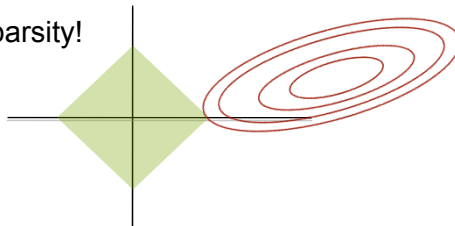$$\hat{\boldsymbol{\beta}} = \mathrm{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2$$

subject to:

$$\sum_{j=1}^{p} |\beta_j| \leq C$$

**Lagrangian Form**

$$\hat{\boldsymbol{\beta}} = \mathrm{argmin}_{\beta} \|\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_1$$

- Still enforces sparsity!
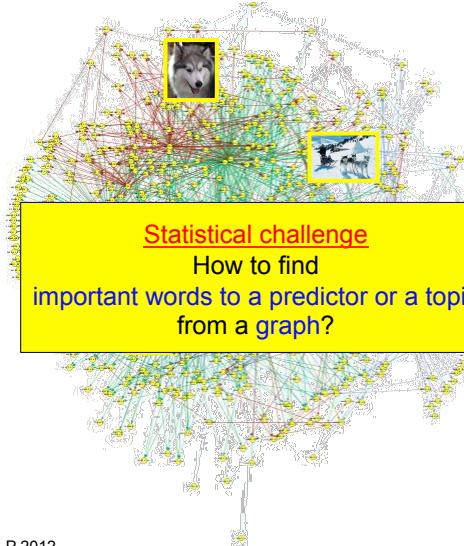
Lasso for Sparse Regression

Class Label    Word counts    Feature strength

$$\beta^* = \arg\min_{\beta}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) \; + \; \lambda\sum_{j=1}^{J}|\beta_j|$$

Lasso Penalty for sparsity

Many zero associations (sparse results), but what if the problem has "structures"?

EMNLP 2012

17

---



Input Structure: the WordNet

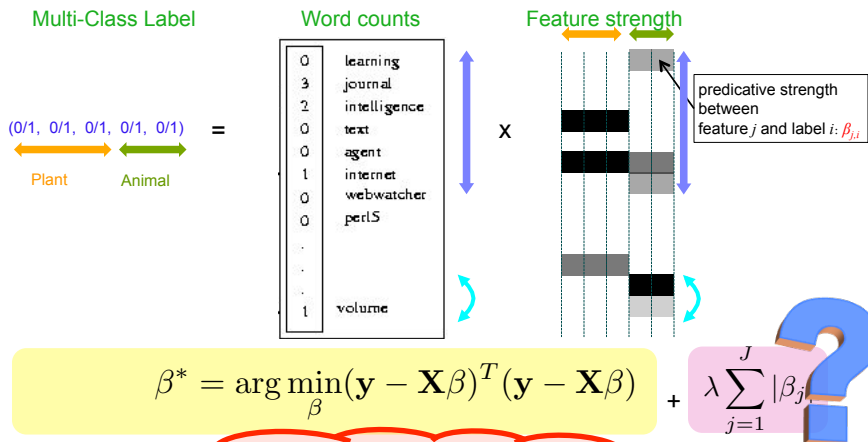- The "network" of synsets in NL
  - Nodes (synsets) represent distinct concept
  - Links represent conceptual-semantic and lexical relations
  - Hidden knowledge and structure among concepts

  - Prior knowledge
  - Context

Statistical challenge
How to find
important words to a predictor or a topic
from a graph?

EMNLP 2012

18

9

# Output Structure: Task Hierarchy

- E.g., the tree hierarchy in the DMOZ repository of the PASCAL Large Scale Hierarchical Text Classification challenge
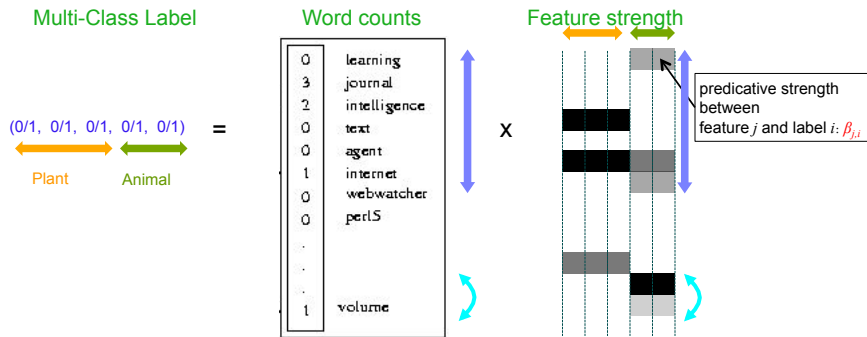


**Statistical challenge**
How to train
multiple labelers that are "related"
by a tree?

---

# Sparse Structured Input/Output Lasso for Multi-task Learning



Multi-Class Label            Word counts            Feature strength

(0/1, 0/1, 0/1, 0/1, 0/1)   =

Plant      Animal

| | |
|---|---|
| 0 | learning |
| 3 | journal |
| 2 | intelligence |
| 0 | text |
| 0 | agent |
| 1 | internet |
| 0 | webwatcher |
| 0 | perlS |
| . | |
| . | |
| . | |
| 1 | volume |

X

predicative strength between feature $j$ and label $i$: $\beta_{j,i}$

$$\beta^* = \arg\min_\beta (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) \; + \; \lambda \sum_{j=1}^{J} |\beta_j|$$

How to combine information across multiple features/classes to increase the power?

10

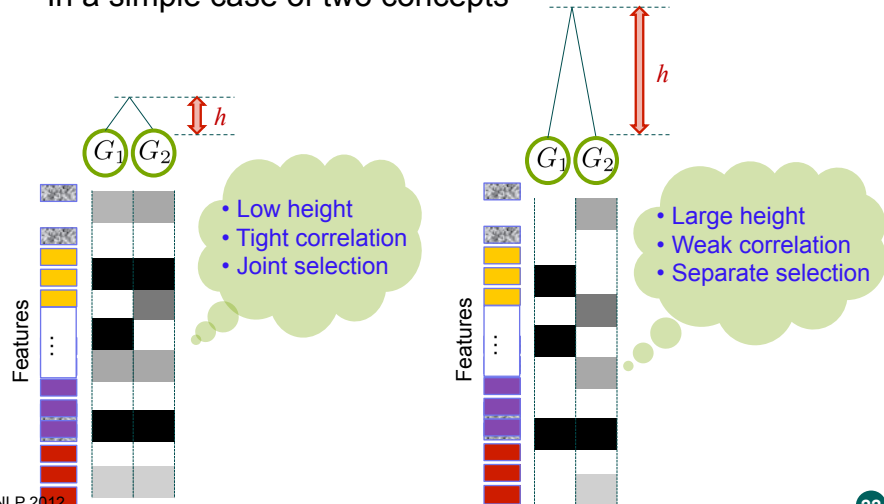## Sparse Structured Input/Output Lasso for Multi-task Learning

Multi-Class Label    Word counts    Feature strength

(0/1, 0/1, 0/1, 0/1, 0/1) =

Plant    Animal

| 0 | learning |
| 3 | journal |
| 2 | intelligence |
| 0 | text |
| 0 | agent |
| 1 | internet |
| 0 | webwatcher |
| 0 | perlS |
| . | . |
| . | . |
| . | . |
| 1 | volume |

X

predicative strength between feature $j$ and label $i$: $\beta_{j,i}$

$$\beta^* = \arg\min_{\beta}(\mathbf{y} - \mathbf{X}\beta)^T(\mathbf{y} - \mathbf{X}\beta) \quad + \quad \lambda\sum_{j=1}^{J}|\beta_j|$$

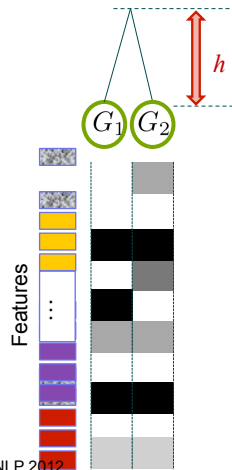+ We introduce
Structured fusion and/or group norm penalties

EMNLP 2012

21

---



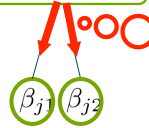# Tree-Guided Group Lasso

• In a simple case of two concepts

$h$

$G_1$ $G_2$

Features

• Low height
• Tight correlation
• Joint selection

$h$

$G_1$ $G_2$

Features

• Large height
• Weak correlation
• Separate selection

EMNLP 2012

[Kim and Xing, ICML 2009]    22

# Tree-Guided Group Lasso

- In a simple case of two concepts

$C_1 = \{\beta_{j1}, \beta_{j2}\}$

$h$

$G_1$ $G_2$

$\beta_{j1}$ $\beta_{j2}$

Select the child nodes jointly or separately?

Features

**Tree-guided group lasso**

$$\text{argmin } (y - X\beta)' \cdot (y - X\beta)$$

$$+ \lambda \sum_j \left[ h\big(|\beta_{j1}| + |\beta_{j2}|\big) + (1 - h)\big(\sqrt{\beta_{j1}^2 + \beta_{j2}^2}\big) \right]$$
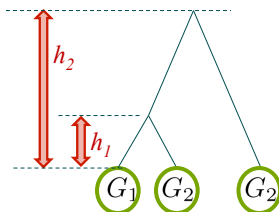
$L_1$ penalty
- Lasso penalty
- Separate selection

$L_2$ penalty
- Group lasso
- Joint selection

Elastic net

EMNLP 2012

**23**

---



# Tree-Guided Group Lasso

- For a general tree

$h_2$

$h_1$

$C_2 = \{\beta_{j1}, \beta_{j2}, \beta_{j3}\}$

$C_1 = \{\beta_{j1}, \beta_{j2}\}$

Select the child nodes jointly or separately?

$G_1$ $G_2$ $G_2$

$\beta_{j1}$ $\beta_{j2}$ $\beta_{j3}$

**Tree-guided group lasso**

$$\text{argmin } (y - X\beta)' \cdot (y - X\beta)$$

$$+ \lambda \sum_j \left[ (1 - h_2)\big(\sqrt{\beta_{j1}^2 + \beta_{j2}^2 + \beta_{j3}^2}\big) + h_2\big(|C_1| + |\beta_{j3}|\big) \right]$$
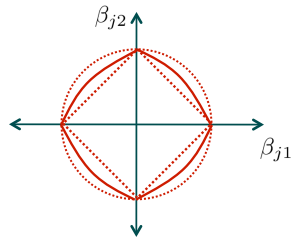
$$(1 - h_1)\big(\sqrt{\beta_{j1}^2 + \beta_{j2}^2}\big) + h_1\big(|\beta_{j1}| + |\beta_{j2}|\big)$$

Joint selection

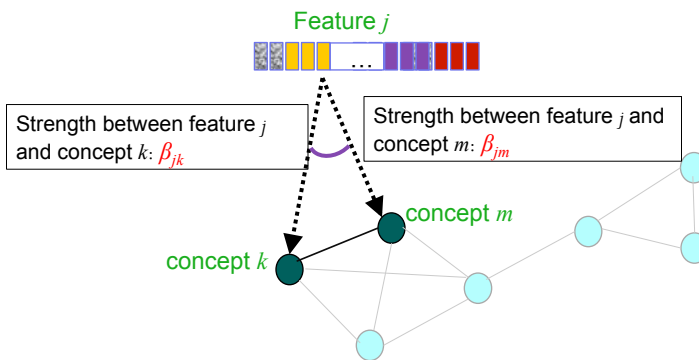Separate selection

EMNLP 2012

**24**

12

**Proposition 1** *For each of the k-th output (gene), the sum of the weights $w_v$ for all nodes $v \in V$ in $T$ whose group $G_v$ contains the k-th output (gene) as a member equals one. In other words, the following holds:*

$$\sum_{v:k \in G_v} w_v = \prod_{m \in Ancestors(v_k)} h_m + \sum_{l \in Ancestors(v_k)} (1 - h_l) \prod_{m \in Ancestors(v_l)} h_m = 1.$$



Previously, in Jenatton, Audibert & Bach, 2009

EMNLP 2012

25

---

# Graph-Guided Fused Lasso

Feature $j$

Strength between feature $j$ and concept $k$: $\beta_{jk}$

Strength between feature $j$ and concept $m$: $\beta_{jm}$
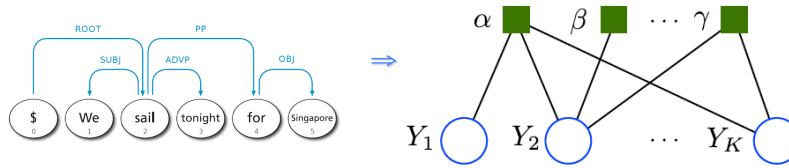
concept $m$

concept $k$



- Fusion Penalty: $|\beta_{jk} - \beta_{jm}|$
- For two correlated concepts (connected in the network), the association strengths may have similar values.
  - Fusion effect propagates to the entire network
  - Association between features and subnetworks of concepts

EMNLP 2012

26

[Kim and Xing, PLOS G 2009]

13

## Full GM-based Loss Functions

$$y^\star = h(\mathbf{x}) \triangleq \arg\max_{y \in \mathcal{Y}} F(\mathbf{x}, y; \mathbf{w})$$
$$F(\mathbf{x}, y; \mathbf{w}) = g(\mathbf{w}^\top \mathbf{f}(\mathbf{x}, y))$$

Represent **factorization** assumptions: $P(\mathbf{y}|x) = \frac{1}{Z} \prod_i \Psi_i(y_i) \prod_\alpha \Psi_\alpha(\mathbf{y}_\alpha)$

**Inference:** compute the MAP, marginals $\mu_i(y_i)$ and $\mu_\alpha(\mathbf{y}_\alpha)$, $Z$

- tractable when $\mathscr{G}$ is a tree, often intractable otherwise

27

---

## Optimization

Original Problem:
$$\arg\max_\beta \equiv \mathcal{L}(\{\mathbf{x}_i, \mathbf{y}_i\}; \beta) + \Omega(\beta)$$
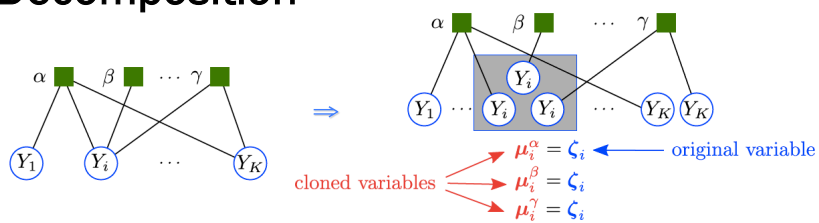
Existing Methods:

| | | |
|---|---|---|
| Interior-point Method (IPM) for Second-order Cone Programming (SOCP) or Quadratic Programming (QP) | 2nd-order, computationally heavy | $\lambda \sum_{g \in \mathcal{G}} w_g \|\boldsymbol{\beta}_g\|_2 \Longrightarrow \lambda \sum_{g \in \mathcal{G}} w_g t_g$  s.t. $\|\boldsymbol{\beta}_g\| \le t_g$ |
| Block Coordinate Descent | Cannot be easily be applied. Hard to compute the subgradient | Optimize $\boldsymbol{\beta}_g$ at one time |

28

14

# New Optimization Framework

- Main Difficulties:
  - Complex loss $\mathcal{L}(\{\mathbf{x}_i, \mathbf{y}_i\}; \beta)$, (e.g., GMs with intractable factors or loopy graphs)
    - Intractable inference
  - Complex shrinkage $\Omega(\beta)$, (e.g., overlapping group penalties)
    - Non-differentiable, non-separatable
- Our approaches:
  - Alternating Direction Dual Decomposition (AD³) [Martins et al, ICML 2011]
  - Proximal Gradient [Chen et al, AOAS 2012]
  - Hierarchical Group Threshholding [Lee and Xing, 2012, submitted]
- Large number of training examples
  - Parallel computation
  - Map-Reduce on computing gradient
    - Map: calculate gradient on single example
    - Reduce: gather gradients computed by all map procedures, and calculate the sum
  - New multi-core framework ..

---

# Alternating Directions Dual Decomposition

[Martins, Figueiredo, Aguiar, Smith, Xing, ICML 2011]



$\mu_i^\alpha = \zeta_i$ ← original variable

cloned variables → $\mu_i^\beta = \zeta_i$

$\mu_i^\gamma = \zeta_i$

$$\text{maximize} \quad \underbrace{\boldsymbol{\theta}_\alpha \cdot \boldsymbol{\mu}_\alpha + \sum_{i \in \mathcal{N}(\alpha)} \left( \frac{\boldsymbol{\theta}_i}{\mathcal{N}(i)} + \boldsymbol{\lambda}_i^\alpha \right) \cdot \boldsymbol{\mu}_i^\alpha}_{\text{linear term}} - \underbrace{\frac{\eta}{2} \sum_{i \in \mathcal{N}(\alpha)} \|\boldsymbol{\mu}_i^\alpha - \boldsymbol{\zeta}_i^\alpha\|^2}_{\text{penalty term}}$$

$$\text{w.r.t.} \quad \boldsymbol{\mu}|_\alpha \in \mathbf{MARG}(\mathscr{G}|_\alpha)$$

- Convergent to the primal and dual solutions (Glowinski and Le Tallec, 1989)
- $O(1/\varepsilon)$ iterations suce for $\varepsilon$-accurate objective (He and Yuan, 2011)
- Solution is always sparse (only $O(|N(\alpha)|)$ nonzeros)
- Active set methods: seek the support of the solution by adding/removing components; very suitable for warm-starting (Nocedal and Wright, 1999)

# Smooth Proximal Gradient Descent

[Chen et al and Xing, UAI 2011, AOAS 2012]

**Original Problem:**

$$\arg\min_{\boldsymbol{\beta}\in\mathbb{R}^J} f(\boldsymbol{\beta}) \equiv \frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \Omega(\boldsymbol{\beta})$$

$$\Omega(\boldsymbol{\beta}) = \max_{\boldsymbol{\alpha}\in\mathcal{Q}} \boldsymbol{\alpha}^T C \boldsymbol{\beta}$$

Separating overlapping constraints

**Approximation Problem:**

$$\arg\min_{\boldsymbol{\beta}\in\mathbb{R}^J} \widetilde{f}(\boldsymbol{\beta}) \equiv \frac{1}{2}\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + f_\mu(\boldsymbol{\beta})$$

$$f_\mu(\boldsymbol{\beta}) = \max_{\boldsymbol{\alpha}\in\mathcal{Q}} \boldsymbol{\alpha}^T C \boldsymbol{\beta} - \mu d(\boldsymbol{\alpha})$$

Smoothing non-differentiable objective

**Gradient of the Approximation:**

$$\nabla \widetilde{f}(\boldsymbol{\beta}) = \mathbf{X}^T(\mathbf{X}\boldsymbol{\beta} - \mathbf{y}) + C^T \boldsymbol{\alpha}^*$$

$$\boldsymbol{\alpha}^* = \arg\max_{\boldsymbol{\alpha}\in\mathcal{Q}} \boldsymbol{\alpha}^T C \boldsymbol{\beta} - \mu d(\boldsymbol{\alpha})$$

$\nabla\widetilde{f}(\boldsymbol{\beta})$ is Lipschitz continuous with the Lipschitz constant $L$
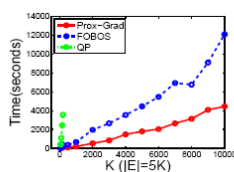
$$L = \lambda_{\max}(\mathbf{X}^T\mathbf{X}) + L_\mu$$

EMNLP 2012

**31**

---

# Convergence Rate

**Theorem**: If we require $f(\boldsymbol{\beta}^t) - f(\boldsymbol{\beta}^*) \leq \epsilon$ and set $\mu = \frac{\epsilon}{2D}$, the number of iterations is upper bounded by:
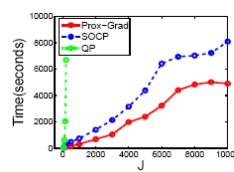
$$t \leq \sqrt{\frac{4\|\boldsymbol{\beta}^*\|_2^2}{\epsilon}\left(\lambda_{\max}(\mathbf{X}^T\mathbf{X}) + \frac{2D\|\Gamma\|^2}{\epsilon}\right)} = O(\frac{1}{\epsilon})$$

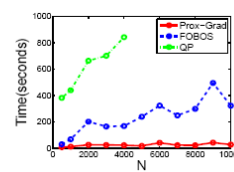Remarks: state of the art IPM method for for SOCP converges at a rate $O(\frac{1}{\epsilon^2})$

**Time complexity** (Per-iteration): $O(J^2K + J\sum_{g\in\mathcal{G}}|g|)$ vs. $O\left(J^2(K+|\mathcal{G}|)^2(KN + J(\sum_{g\in\mathcal{G}}|g|))\right)$



(a)   (b)   (c)

EMNLP 2012

**32**

16

# What if the structure becomes too complex ?

- Too many groups in real problems:

$$\beta_{io-lasso} = \arg\min_{\beta} \sum_{k=1}^{K} \sum_{i=1}^{N} \left( Y_i^k - \sum_{j=1}^{p} \beta_j^k X_{ij} - \sum_{(r,s)\in U} \beta_{rs}^k Z_{i,rs} \right) + \lambda_1 \sum_{j} \sum_{k=1} \left| \beta_j^k \right|$$

$$+ \lambda_2 \sum_{k} \sum_{m} \sqrt{\sum_{(r,s)\in S_m} \beta_{rs}^{k^2}}$$

$$+ \lambda_3 \sum_{j} \sqrt{\sum_{k} \beta_j^{k^2}}$$

$$+ \lambda_4 \sum_{k} \sum_{(r,s)\in I} \left| \beta_{rs}^k \right|$$

Output structure, group selection of
lasso penalty with group sparsity
features across multiple predictors
Easily hundreds or thousands of group constraints!

- Recall that even SPG has a complexity of $O(J^2 K + J \sum_{g \in \mathcal{G}} |g|)$

- And an optimization procedure must
  - minimize our objective function, and
  - induce correct sparsity patterns

- Hierarchical group-thresholding:
  - an algorithmic approach to directly reduce search space of sparsity, while optimizes the exact loss
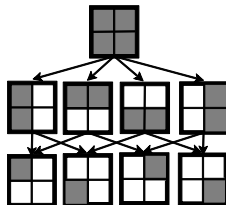
33

---



# Hierarchical Group-Thresholding

[Lee and Xing, Submitted 2012]

- DAG for Sparsity Patterns
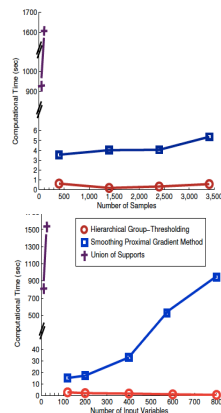  - All sparsity patterns of a 2x2 matrix:
  - A DAG of inclusion relation relationships of sparsity

- Hierarchical Group-Thresholding
  - Initialize B using ridge regression
  - Step 1: Traversing DAG, check the optimality condition of the zero pattern at each node. If the condition holds, set zero
  - Step 2: Update non-zero regression coefficients using coordinate descent

34

17

# What if the function is non-linear?

**Group Sparse Additive Models** [Ying, Chen, Xing, ICML 2012]

- Assume G is a partition of {1, · · · , p}, i.e., the groups in G do not overlap.
- The optimization problem is

- $$\min_{\mathbf{f}} L(\mathbf{f}) + \lambda \Omega_{\mathrm{group}}(\mathbf{f}),$$

  where
  $$\Omega_{\mathrm{group}}(\mathbf{f}) = \sum_{g \in \mathcal{G}} \sqrt{|g|} \|\mathbf{f}_g\| = \sum_{g \in \mathcal{G}} \sqrt{|g|} \sqrt{\sum_{j \in g} \mathbb{E}\left[f_j^2(X_j)\right]}.$$

- Non-trivial to solve due to
  - correlation structure of component functions within the group
  - non-smoothness of functional group penalty

---

# Toward Human-Level Intelligence

- Now we have dealt with high feature dimension
  - Sparsity

- and we have know how to leverage structural knowledge
  - Structured shrinkage

- What about massive **concept space**?

The IM GENET knowledge ontology



mammal → placental → carnivore → canine → dog → working dog → husky

Data: >20 million images

Features: ~1 million (number comes from the top performing system in ILSVRC10, [Lin et al. 2011])

Classes: ~22k classes

# Output Coding

| M | 1 | 2 | 3 | ... | ... | K |
|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| 2 | -1 | 0 | 0 | 1 | 1 | 0 |
| ... | 0 | -1 | 0 | -1 | 0 | 1 |
| C | 0 | 0 | -1 | 0 | -1 | -1 |

- Every class is now represented by a bit-string
  - *Coding*: a codeword is assigned to each class
  - *Decoding*: given test data, look for most similar class codeword
- Predict bit by bit through binary or ternary classifier – this is much easier than the 1 vs C-1 classifier
- Decoding the bit-string – error correcting

EMNLP 2012

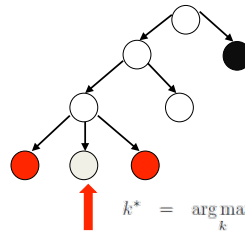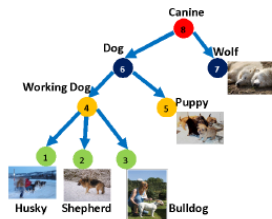[Zhao and Xing, 2012, submitted]

37

---

| M | 1 | 2 | 3 | ... | ... | K |
|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| 2 | -1 | 0 | 0 | 1 | 1 | 0 |
| ... | 0 | -1 | 0 | -1 | 0 | 1 |
| C | 0 | 0 | -1 | 0 | -1 | -1 |

# Learning the Coding Matrix

- Accuracy of base binary classifiers for bit-prediction
  - Use category hierarchy for a measure of separability
  - Large intra-partite similarity + small inter-partite similarity
- Strong error-correcting ability
  - Maximize distance between rows of coding matrix
- Fault tolerance
  - Introduction of ignored classes: {-1,0,+1} instead of {-1,+1}

$$\max_{\mathbf{B}} \quad F_b(\mathbf{B}) - \lambda_r F_r(\mathbf{B}) - \lambda_c \sum_{l=1}^{L} ||\boldsymbol{\beta}_l||_2^2$$

$$s.t. \quad \mathbf{B} \in \{-1, 0, +1\}^{K \times L}$$

$$\sum_{k=1}^{K} (|B_{kl}| + B_{kl}) \geq 2, \ \forall l = 1, \ldots, L$$

$$\sum_{k=1}^{K} (|B_{kl}| - B_{kl}) \geq 2, \ \forall l = 1, \ldots, L$$

$$\sum_{l=1}^{L} |B_{kl}| \geq 1, \ \forall k = 1, \ldots, K$$

EMNLP 2012

38

19

# Probabilistic Decoding

- Output code can have real semantic meaning
  - E.g., encoding a tree path in a label taxonomy
- Probabilistic decoding:
  - bit $i$ depend on bit $j$ probabilistically
- Define prior $P(y_l | \hat{y} = k)$ using tree hierarchy
  - Graph coloring: all nodes participating in $i$-th bit prediction are colored (red for positive, black for negative)
  - Task: what is the probability of node $k$ being colored red?



$$k^\star = \arg\max_k P(\hat{y} = k | w_1, \ldots, w_L, x)$$

EMNLP 2012

39

---

# Multi-Way Classification Accuracy

- DMOZ repository in PASCAL Hierarchical Text Classification challenge

| Data set | # classes | # training data | # test data | # features |
|----------|-----------|-----------------|-------------|------------|
| DMOZ-small | 1139 | 6323 | 1858 | 1199848 |
| DMOZ-large | 12294 | 93805 | 34905 | 1199856 |

| Algorithm | DMOZ small | | DMOZ large | |
|-----------|-------|-------|-------|-------|
|  | Top 1 | Top 5 | Top 1 | Top 5 |
| OVR | 50.91 | 64.72 | 37.24 | 44.87 |
| RDOC | 41.77 | 54.52 | 5.13 | 7.99 |
| RSOC | 42.30 | 58.40 | 5.47 | 7.23 |
| SpectralOC | 44.83 | 59.10 | 22.10 | 23.82 |
| SSOC | **56.67** | **67.33** | **41.28** | **46.71** |

EMNLP 2012

[Zhao and Xing, Submitted 2012]   40

20

## Discovering Sociolinguistic Associations on Twitter

- Twitter Gardenhose feed from March 1-7, 2010
- 9250 authors, 380,000 messages, 4.7 million tokens
- Filters:
  - At least 20 messages (in Gardenhose)
  - Messages must include GPS within a USA zipcode
  - No more than 1000 followers, followees
- GPS → Zipcode → U.S. Census Demographic Statistics
  - Zipcodes commonly proxy for demographics in public health.
  - Careful! Twitter users are not an unbiased sample from a zipcode.

[Eisenstein, Smith, and Xing. ACL 2011]  41

---

## Demographic multi-prediction

$$X \quad \Theta \quad \approx \quad Y$$



$L1/L\infty$ regularizer

$$\Omega(\boldsymbol{\Theta}) = \sum_t \max_p |\theta_{pt}|$$

aka **multi-output lasso**
(Turlach et al 2005)

| vocabulary | # features | average | white | Afr. Am. | Hisp. | Eng. lang. | Span. lang. | other lang. | urban | family | renter | med. inc. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| full | 5418 | 0.260 | 0.337 | 0.318 | 0.296 | 0.384 | 0.296 | 0.256 | 0.155 | 0.113 | 0.295 | 0.152 |
| multi-output lasso | | 0.260 | 0.326 | 0.308 | 0.304 | 0.383 | 0.303 | 0.249 | 0.153 | 0.113 | 0.302 | 0.156 |
| SVD | 394.6 | 0.237 | 0.321 | 0.299 | 0.269 | 0.352 | 0.272 | 0.226 | 0.138 | 0.081 | 0.278 | 0.136 |
| highest variance | | 0.220 | 0.309 | 0.287 | 0.245 | 0.315 | 0.248 | 0.199 | 0.132 | 0.085 | 0.250 | 0.135 |
| most frequent | | 0.204 | 0.294 | 0.264 | 0.222 | 0.293 | 0.229 | 0.178 | 0.129 | 0.073 | 0.228 | 0.126 |

[Eisenstein, Smith, and Xing. ACL 2011]  42

21

# Sociolinguistic Associations

| | white | Afr. Am. | Hisp. | Eng. lang. | Span. lang. | other lang. | urban | family | renter | med. inc. |
|---|---|---|---|---|---|---|---|---|---|---|
| as | | | - | + | - | | | | | |
| awesome | + | - | | | | | - | | - | + |
| break | | | - | + | - | - | | | | |
| campus | | | - | + | - | - | | | | |
| dead | - | + | | - | + | | + | | + | |
| hell | | - | | + | - | - | | | | |
| shit | - | | | | | | | | + | |
| train | | | | - | + | | | | + | |
| will | | | - | + | - | | | | | |
| would | | | | + | | | | | | - |

| | white | Afr. Am. | Hisp. | Eng. lang. | Span. lang. | other lang. | urban | family | renter | med. inc. |
|---|---|---|---|---|---|---|---|---|---|---|
| -_- | - | | + | - | + | + | + | | | |
| ;) | | - | + | - | + | | | | | |
| :( | | - | | | | | | | | |
| :) | | - | | | | | | | | |
| :d | + | - | + | - | + | | | | | |

| | white | Afr. Am. | Hisp. | Eng. lang. | Span. lang. | other lang. | urban | family | renter | med. inc. |
|---|---|---|---|---|---|---|---|---|---|---|
| bbm | - | + | | - | | | + | + | | + |
| lls | | + | - | + | - | - | | | | |
| lmaoo | - | + | + | - | + | + | + | | + | |
| lmaooo | - | + | + | - | + | + | + | | + | |
| lmaoooo | - | + | + | - | + | + | | | + | |
| lmfaoo | - | | + | - | + | + | | | + | |
| lmfaooo | - | | + | - | + | + | | | + | |
| lml | - | + | + | - | + | + | + | | + | - |
| odee | - | | + | - | + | | + | | + | |
| omw | - | + | + | - | + | + | + | | + | |
| smfh | - | + | + | - | + | + | + | | + | |
| smh | - | + | | | | | | + | + | |
| w| | - | | + | - | + | + | + | | + | |

[Eisenstein, Smith, and Xing. ACL 2011] 43

---

# Dependency Parsing

Datasets from CoNLL-2006 and CoNLL-2008 shared tasks



Legend: Best @CONLL-X, Best, Our Model. UAS (%) by language: Arabic, Bulgarian, Chinese, Czech, English, Danish, Dutch, German, Japanese, Portuguese, Slovene, Spanish, Swedish, Turkish, PTB

- ■ Best published includes:
  - ■ **transition-based** (Nivre et al., 2006; Huang and Sagae, 2010),
  - ■ **graph-based** (McDonald et al., 2006; Koo and Collins, 2010),
  - ■ **hybrid** (Nivre and McDonald, 2008; Martins et al., 2008a),
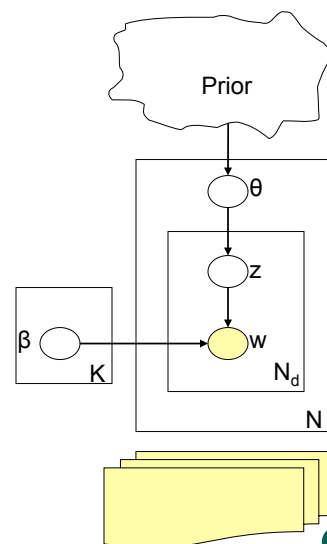  - ■ **turbo parsers** (Martins et al., 2010; Koo et al., 2010)

[Martins, Smith, and Xing. ACL 2009, Martins, EMNLP 2010, 2011] 44

# Outline

- Sparse Structured Input-Output Models
  - … supervised learning
  - … convex optimization and log loss
  - … Frequentist-style shrinkage via regularization

- Sparse Topic Models
  - … unsupervised learning
  - … non-convex and likelihood-driven
  - … Bayesian-style posterior inference

- Sparse and Discriminative Topic Models?
  - … toward jointly explorative and predictive learning

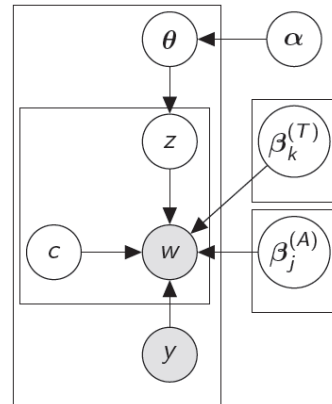# Modeling Semantics: e.g., Topic Models

Generating a document

> – *Draw $\theta$ from the prior*
>
> For each word $n$
>
> - Draw $z_n$ from $multinomial(\theta)$
>
> - Draw $w_n \mid z_n, \{\beta_{1:k}\}$ from $multinomial(\beta_{z_n})$

- **Prior over topic Vector**
  - Latent Dirichlet Allocation (LDA)
  - Correlated priors (CTM)
  - Hierarchical priors
- **Topics**
  - Unigram, bigrams, etc
- **Document structure**
  - Bag of words
  - Multi-modal



Side information

## Modeling and Inference Complexity

- What if we want to combine latent topics with additional facets, such as geography in a unsupervised fashion?

- Additional latent variables decide which facet is responsible for each token (e.g. Ahmed and Xing 2010).

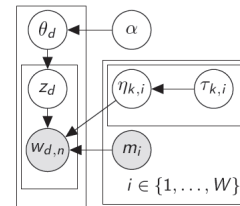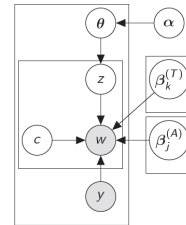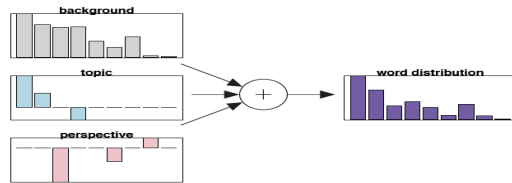- That's twice as many latent variables per document!

---

## Compact modes needed on mobile devices

Sparse Additive Generative Models

EMNLP 2012
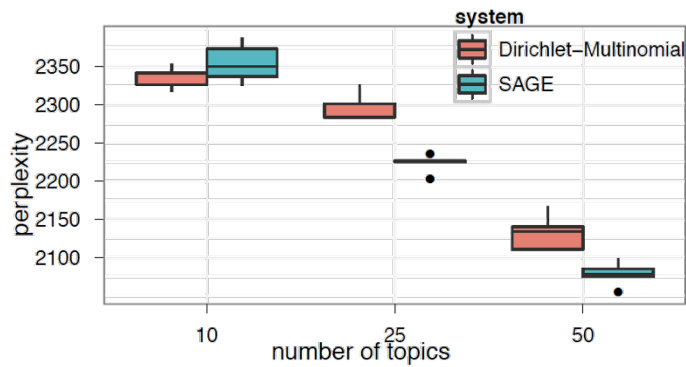
J. Eisenstein, A. Ahmed, E.P. Xing. ICML, 2011 · 49



# Model Compression on Text

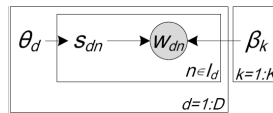- NIPS dataset: 1986 training docs, 10K vocabulary

- Adaptive sparsity:
  - 5% non-zeros for 10 topics
  - 1% non-zeros for 50 topics

EMNLP 2012

50

# Sparse Topical Coding

- Goal: design a non-probabilistic topic model that is amenable to
  - direct control on the posterior sparsity of inferred representations
  - avoid dealing with normalization constant when considering supervision or rich features
  - seamless integration with a convex loss function (e.g., svm hinge loss)

- We extend sparse coding to hierarchical sparse topical coding
  - word code $\theta$
  - document code $\boldsymbol{s}$

$$\theta_d \rightarrow s_{dn} \rightarrow w_{dn} \leftarrow \beta_k$$
$$n \in I_d \quad k=1{:}K$$
$$d=1{:}D$$

reconstruction loss           sparse codes

$$\min_{\{\boldsymbol{\theta}_d,\mathbf{s}_d\},\boldsymbol{\beta}} \sum_{d,n\in I_d} \ell(w_{dn},\mathbf{s}_{dn}^\top \boldsymbol{\beta}_{\cdot n}) + \lambda \sum_d \|\boldsymbol{\theta}_d\|_1 + \sum_{d,n\in I_d} (\gamma\|\mathbf{s}_{dn}-\boldsymbol{\theta}_d\|_2^2 + \rho\|\mathbf{s}_{dn}\|_1)$$

$$\text{s.t.} : \boldsymbol{\theta}_d \geq 0, \ \mathbf{s}_{dn} \geq 0, \ \forall d,n \in I_d; \ \boldsymbol{\beta}_k \in \mathcal{P}, \ \forall k,$$

truncated aggregation

non-negative codes       topical bases

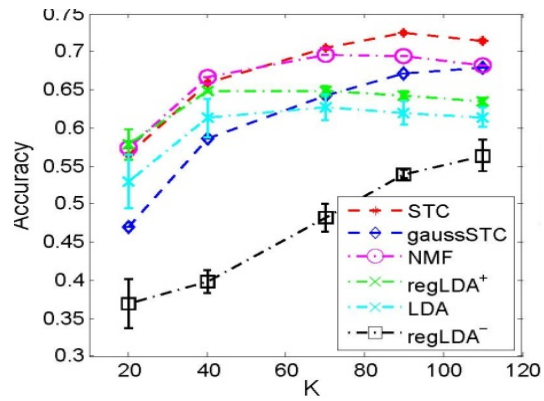J. Zhu, & E.P. Xing. UAI, 2011   51

---

# Algorithms on Sparse Latent Space Models

- Complex objective
  - Non-convex, but often bi-convex
  - Often additional non-negativity constraints other than sparsity

- Hierarchical sparse coding
  - Greedy algorithm for the non-convex $L_0$ "pseudo-norm":
    - select the element with maximum correlation with the residual
    - known as "matching pursuit" (Mallat & Zhang, 1993)
  - For the convex $L_1$ norm, many algorithms:
    - Soft-thresholding with coordinate descent (Friedman et al., 2007; Zhu & Xing, 2011)
    - Proximal methods (Nesterov, 2007; Jenatton et al., 2010, Chen et al 2011)
    - Active-set methods (Roth & Fischer, 2008)
    - Online/stochastic variants
    - …

- Dictionary (topic) learning
  - projected gradient descent
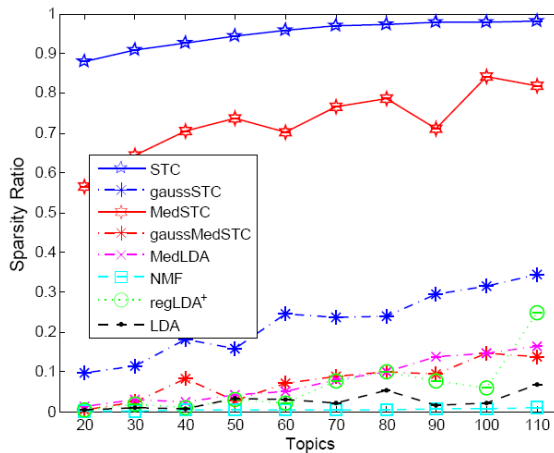  - any faster alternative method can be used

52

# Comparisons

LDA vs. STC
[Zhu and Xing, UAI 2011]

53

# Sparse word codes

- Sparsity ratio: percentage of zeros



- NMF: non-negative matrix factorization
- MedLDA (Zhu et al., 2009)
- regLDA: LDA with entropic regularizer
- gaussSTC: use L2 rather than L1-norm

54

# Outline

- Sparse Structured Input-Output Models
  - … supervised learning
  - … convex optimization and log loss
  - … Frequentist-style shrinkage via regularization

- Sparse Topic Models
  - … unsupervised learning
  - … non-convex and likelihood-driven
  - … Bayesian-style posterior inference

- Sparse and Discriminative Topic Models?
  - … toward jointly explorative and predictive learning

---

## *Predictive* Subspace Learning with *Supervision*

- Unsupervised latent subspace representations are generic but can be sub-optimal for predictions
- Many datasets are available with supervised side information
  - Tripadvisor Hotel Review ( http://www.tripadvisor.com )
  - LabelMe http://labelme.csail.mit.edu/
    - Many others
  - Flickr (http://www.flickr.com/)
  - IMAGENET
- Can be noisy, but not random noise (Ames & Naaman, 2007)
  - labels & rating scores are usually assigned based on some intrinsic property of the data
  - helpful to suppress noise and capture the most useful aspects of the data
- **Goals:**
  - **Discover latent subspace representations that are both *predictive* and *interpretable* by exploring weak supervision information**

28

# MLE versus Max-margin Learning

- Likelihood-based estimation
  - Probabilistic (joint/conditional likelihood model)
  - Easy to perform Bayesian learning, and incorporate prior knowledge, latent structures, missing data
  - Bayesian or direct regularization
  - Hidden structures or generative hierarchy

- Max-margin learning
  - Non-probabilistic (concentrate on input-output mapping)
  - Not obvious how to perform Bayesian learning or consider prior, and missing data
  - Support vector property, sound theoretical guarantee with limited samples
  - Kernel tricks

- Maximum Entropy Discrimination (MED) (Jaakkola, et al., 1999)
  - Model averaging $\hat{y} = \text{sign} \int p(\mathbf{w}) F(x; \mathbf{w}) \, d\mathbf{w}$ $\quad (y \in \{+1, -1\})$
  - The optimization problem (binary classification)

  $$\min_{p(\Theta)} KL(p(\Theta) || p_0(\Theta))$$

  $$\text{s.t.} \quad \int p(\Theta)[y_i F(x; \mathbf{w}) - \xi_i] \, d\Theta \geq 0, \forall i,$$

  *where $\Theta$ is the parameter $\mathbf{w}$ when $\xi$ are kept fixed or the pair $(\mathbf{w}, \xi)$ when we want to optimize over $\xi$*

57

---

# MaxEnt Discrimination Markov Network

(Zhu et al, ICML 2008, Zhu and Xing, JMLR 2009)

- **Structured MaxEnt Discrimination (SMED):**

  $$\text{P1}: \quad \min_{p(\mathbf{w}), \xi} \boxed{KL(p(\mathbf{w}) || p_0(\mathbf{w})) + U(\xi)}$$

  $$\text{s.t.} \quad p(\mathbf{w}) \in \mathcal{F}_1, \ \xi_i \geq 0, \forall i.$$

  *generalized* maximum entropy or *regularized* KL-divergence

- **Feasible subspace of weight distribution:**

  $$\mathcal{F}_1 = \left\{ p(\mathbf{w}) : \boxed{\int p(\mathbf{w})[\Delta F_i(\mathbf{y}; \mathbf{w}) - \Delta \ell_i(\mathbf{y})] \, d\mathbf{w} \geq -\xi_i,} \forall i, \forall \mathbf{y} \neq \mathbf{y}^i \right\},$$

  *expected* margin constraints.

- **Average from a distribution of M³Ns**

  $$h_1(\mathbf{x}; p(\mathbf{w})) = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \int p(\mathbf{w}) F(\mathbf{x}, \mathbf{y}; \mathbf{w}) dw$$

$p_0$

$D(p, p_0) = KL(p || p_0)$

$p$

58

29

# Maximum Entropy Discrimination LDA (MedLDA)

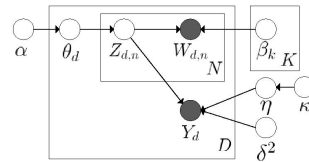- Bayesian sLDA:



- MED Estimation:
  - MedLDA Regression Model

$$P1(\text{MedLDA}^r): \min_{q,\alpha,\beta,\delta^2,\xi,\xi^\star} \mathcal{L}(q) + C\sum_{d=1}^{D}(\xi_d + \xi_d^\star)$$

**model fitting**

**predictive accuracy**

$$\text{s.t. } \forall d: \begin{cases} y_d - E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d, & \mu_d \\ -y_d + E[\eta^\top \bar{Z}_d] \leq \epsilon + \xi_d^\star, & \mu_d^\star \\ \xi_d \geq 0, & v_d \\ \xi_d^\star \geq 0, & v_d^\star \end{cases}$$

  - MedLDA Classification Model

$$P2(\text{MedLDA}^c): \min_{q,q(\eta),\alpha,\beta,\xi} \mathcal{L}(q) + C\sum_{d=1}^{D}\xi_d$$

$$\text{s.t. } \forall d, \ y \neq y_d: \quad E[\eta^\top \Delta\mathbf{f}_d(y)] \geq 1 - \xi_d; \ \xi_d \geq 0.$$

(Zhu et al, ICML 2009, JMLR 2012) **59**

---

# Document Modeling

- Data Set: 20 Newsgroups
- 110 topics + 2D embedding with t-SNE (var der Maaten & Hinton, 2008)



MedLDA                              LDA

**60**

30

# Document Modeling

comp.graphics

politics.mideast

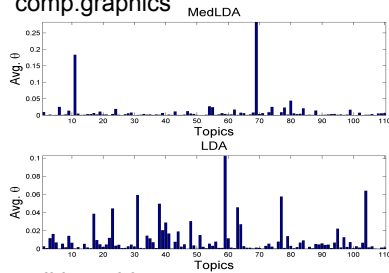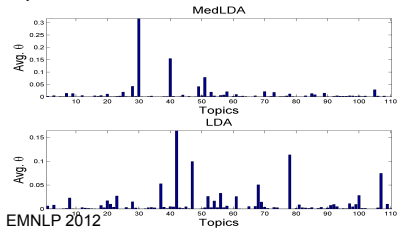| MedLDA | | | LDA | | |
|---|---|---|---|---|---|
| T 69 | T 11 | T 80 | T 59 | T 104 | T 31 |
| image | graphics | db | image | ftp | card |
| jpeg | image | key | jpeg | pub | monitor |
| gif | data | chip | color | graphics | dos |
| file | ftp | encryption | file | mail | video |
| color | software | clipper | gif | version | apple |
| files | pub | system | images | tar | windows |
| bit | mail | government | format | file | drivers |
| images | package | keys | bit | information | vga |
| format | fax | law | files | send | cards |
| program | images | escrow | display | server | graphics |

| T 30 | T 40 | T 51 | T 42 | T 78 | T 47 |
|---|---|---|---|---|---|
| israel | turkish | israel | israel | jews | armenian |
| israeli | armenian | lebanese | israeli | jewish | turkish |
| jews | armenians | israeli | peace | israel | armenians |
| arab | armenia | lebanon | writes | israeli | armenia |
| writes | people | people | article | arab | turks |
| people | turks | attacks | arab | people | genocide |
| article | greek | soldiers | war | arabs | russian |
| jewish | turkey | villages | lebanese | center | soviet |
| state | government | peace | lebanon | jew | people |
| rights | soviet | writes | people | nazi | muslim |

---

# Classification

- **Data Set:** 20Newsgroups
  - Binary classification: "alt.atheism" and "talk.religion.misc" (Simon et al., 2008)
  - Multiclass Classification: all the 20 categories
- **Models**: DiscLDA, sLDA (Binary ONLY! Classification sLDA (Wang et al., 2009)), LDA+SVM (baseline), MedLDA, MedLDA+SVM
- **Measure**: Relative Improvement Ratio

$$RR(\mathcal{M}) = \frac{precision(\mathcal{M})}{precision(LDA + SVM)} - 1$$

31

# Regression

- **Data Set**: Movie Review (Blei & McAuliffe, 2007)
- **Models**: MedLDA(*partial*), MedLDA(*full*), sLDA, LDA+SVR
- **Measure**: predictive R$^2$ and per-word log-likelihood

$$pR^2 = 1 - \frac{\sum_d (y_d - \hat{y}_d)^2}{\sum_d (y_d - \bar{y}_d)^2}$$

---

# Time Efficiency

- Binary Classification



- Multiclass:
  - MedLDA is comparable with LDA+SVM
- Regression:
  - MedLDA is comparable with sLDA

# Supervised STC



- Joint loss minimization

$$\min_{\{\theta_d\},\{\mathbf{s}_d\},\beta,\eta} \quad f(\{\theta_d\},\{\mathbf{s}_d\},\beta) + C\mathcal{R}_h(\{\theta_d\},\eta) + \frac{1}{2}\|\eta\|_2^2$$

$$\text{s.t.:} \quad \theta_d \geq 0,\ \forall d;\ \mathbf{s}_{dn} \geq 0,\ \forall d, n \in I_d;\ \beta_k \in \mathcal{P},\ \forall k,$$

  - coordinate descent alg. applies with closed-form update rules
  - No sum-exp function; seamless integration with non-probabilistic large-margin principle

(Zhu and Xing, UAI 2011)

**65**

---

# Classification accuracy

- 20 newsgroup data:

**66**

## Summary: Margin-based Learning Paradigms

SVM
$$y = \text{sign}(\mathbf{w}^\top \mathbf{x} + b)$$

$$\min_{\mathbf{w},\xi} \frac{1}{2}\mathbf{w}^\top\mathbf{w} + C\sum_{i=1}^{N}\xi_i;$$
$$\text{s.t.} \ \ y_i(\mathbf{w}^\top\mathbf{x}_i + b) \geq 1 - \xi_i, \ \forall i.$$

Structured prediction →

$M^3N$
$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ & \vdots & \end{pmatrix}$$
$$\mathbf{y}^\star = \arg\max_{\mathbf{y}} \mathbf{w}^\top \mathbf{f}(\mathbf{x},\mathbf{y};\mathbf{w})$$
$$\mathbf{y} = \begin{pmatrix} y_{21} & y_{22} & \cdots \\ & \vdots & \cdots \end{pmatrix}$$

$$\min_{\mathbf{w},\xi} \frac{1}{2}\|\mathbf{w}\|^2 + C\sum_{i=1}^{N}\xi_i$$
$$\text{s.t.} : \ \mathbf{w}^\top\Delta\mathbf{f}_i(\mathbf{y}) \geq \Delta\ell_i(\mathbf{y}) - \xi_i, \xi_i \geq 0, \forall i, \forall \mathbf{y} \neq \mathbf{y}^i$$

Bayes learning ↓

MED
$$y = \text{sign}(\langle \mathbf{w}^\top \mathbf{f}(\mathbf{x}) \rangle_{p(\mathbf{w})})$$

$$\min_{p,\xi} KL(p\|p_0) + C\sum_{i=1}^{N}\xi_i;$$
$$\text{s.t.} \ \ y_i\langle \mathbf{f}(\mathbf{x}_i)\rangle_{p(\mathbf{w})} \geq 1 - \xi_i, \ \forall i.$$

Structured prediction →

MaxEnDNet
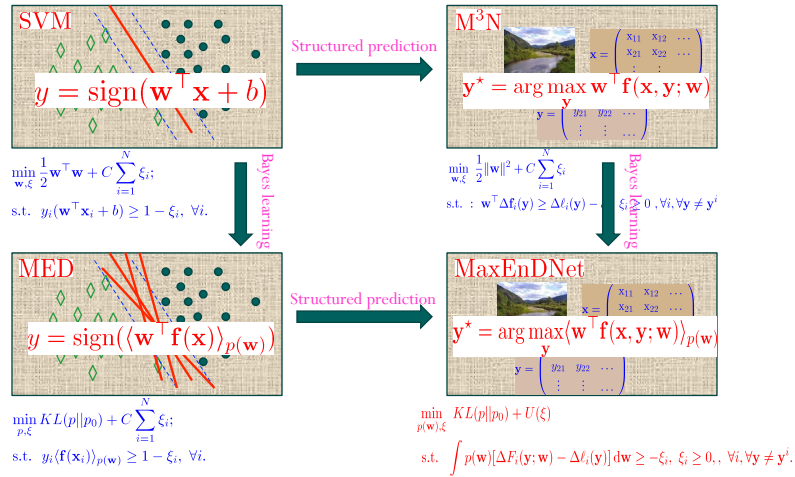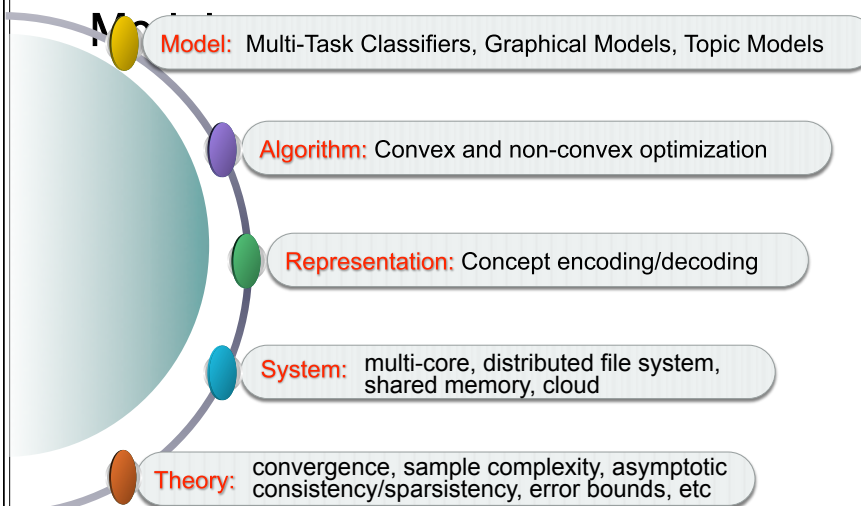$$\mathbf{x} = \begin{pmatrix} x_{11} & x_{12} & \cdots \\ x_{21} & x_{22} & \cdots \\ & \vdots & \end{pmatrix}$$
$$\mathbf{y}^\star = \arg\max_{\mathbf{y}}\langle \mathbf{w}^\top \mathbf{f}(\mathbf{x},\mathbf{y};\mathbf{w})\rangle_{p(\mathbf{w})}$$
$$\mathbf{y} = \begin{pmatrix} y_{21} & y_{22} & \cdots \\ & \vdots & \cdots \end{pmatrix}$$

$$\min_{p(\mathbf{w}),\xi} KL(p\|p_0) + U(\xi)$$
$$\text{s.t.} \ \int p(\mathbf{w})[\Delta F_i(\mathbf{y};\mathbf{w}) - \Delta\ell_i(\mathbf{y})]\,d\mathbf{w} \geq -\xi_i, \ \xi_i \geq 0,, \ \forall i, \forall \mathbf{y} \neq \mathbf{y}^i.$$

Bayes learning ↓

EMNLP 2012

**67**

---



## Conclusion and Challenges: Learning Sparse Structured Input/Output Models

**Model:** Multi-Task Classifiers, Graphical Models, Topic Models

**Algorithm:** Convex and non-convex optimization

**Representation:** Concept encoding/decoding

**System:** multi-core, distributed file system, shared memory, cloud

**Theory:** convergence, sample complexity, asymptotic consistency/sparsistency, error bounds, etc

EMNLP 2012

**68**

# Thanks!

Reference: