Structured Sparsity

Sung Ju Hwang Ulsan National Institute of Science and Technology

Lecturer Introduction

Works on the intersection between Computer Vision and Machine Learning



Expanding Object Detector's Horizon: Incremental Learning Framework for Object Detection in Videos, **CVPR 2015**

Taxonomies, NIPS 2012



A Unified Semantic Embedding: Relating Taxonomies with Attributes, **NIPS 2014**

Features, NIPS 2011



Analogy Preserving Semantic Embedding for Visual Object Categorization, **ICML 2013**





Sharing Features between Objects and Their Attributes, **CVPR 2011**

Lecture Note Update

http://sjhwang.unist.ac.kr/sparsity.pdf

Current Research Interests

Learning Semantics

 How can we learn high-level semantic knowledge from the given visual and textual data, such that machine's understanding of the world aligns well with how we understand the world?

Interactive Learning

- How can we build a learning model such that machines learn from interacting with humans?

Lifelong Learning

- How can we learn a model that can basically learn forever, while transferring knowledge obtained at earlier stages to the learning for later ones?

Part 1: General Sparsity

- Background: Optimization
- Regularizations
- Sparsity
- Sparsity-inducing Regularization
- Example: lasso
- Numerical optimization proximal methods
- Alternative sparse methods
- Example: Image restoration
- Example: Self-taught learning

Background: Mathematical Optimization

Selection of the values that minimize/maximize a given function

minimize $f_0(w)$ subject to $f_i(w) \le b_i, i = 1, ..., m$

$$\min_{w} f_0(w)$$

s.t. $f_i(w) \le b_i, i = 1, ..., m$

 $w = (w_1, ..., w_n)$ optimization variables $f_0: R^n \to R$ objective function $f_i: R^n \to R, i = 1, ..., m$ constraints function

We want to find the optimal solution w^* that has the smallest value of f_0 among all vectors that satisfy the constraints.

Background: Convex Functions

 $f: \mathbb{R}^n \to \mathbb{R}$ is convex if dom f is a convex set and $f(\theta x + (1 - \theta)y) \le \theta f(x) + (1 - \theta)f(y)$



Has a unique global minimum. f is concave if -f is convex

Regularization

Introduction of bias to better condition the target problem with additional information

$$\min_{\boldsymbol{w}} f(\boldsymbol{w}) + \lambda \Omega(\boldsymbol{w})$$



For machine learning, regularization is often used to prevent overfitting.

Regularization

L2 regularizations, based on 2 norm, is one the most common types of regularizations



Shrinks the value of the variables while favoring similar weights among them

Sparsity

We define the LO-psuedonorm as follows:

$$\|\boldsymbol{w}\|_{0} = \#\{i = 1, \dots, p \mid \boldsymbol{w}_{i} \neq 0\}$$

$$w_{1} \quad w_{2} \quad w_{3} \quad w_{4} \quad w_{5} \quad w_{6} \quad w_{7} \quad w_{8}$$

This is a measure of how many of the variables are non-zero. L0-regularization results in learning **sparse** models, by selecting variables

Why is Sparsity Good?

Feature selection - Identifies features that are truly relevant to the target task.



Useful for high-dimensional learning and data-driven methods

Why is Sparsity Good?

Better Interpretability – The learned model can be better explained in terms of selected non-zero entries.

Category	Ground-truth attributes	Supercategory + learned attributes
Otter	An animal that swims, fish, water, new world, small, flippers, furry, black, brown, tail,	A musteline mammal that is quadrapedal, flippers, furry, ocean
Skunk	An animal that is smelly, black, stripes, white, tail, furry, ground, quadrapedal, new world, walks,	A musteline mammal that has stripes
Deer	An animal that is brown, fast, horns, grazer, forest, quadrapedal, vegetation, timid, hooves, walks,	A deer that has spots, nestspot, longneck, yellow, hooves
Moose	An animal that has horns, brown, big, quadrapedal, new world, vegetation, grazer, hooves, strong, ground,	A deer that is arctic, stripes, black
Equine	N/A	An odd-toed ungulate, that is lean and active
Primate	N/A	An animal, that has hands and bipedal

Why is Sparsity Good?

Model compression – With most of the parameters set to zero, sparsity can greatly reduce the memory and computational requirements.



8x8 matrix – requires 512 bytes to store in double precision

12 nonzero entries – requires 96 bytes + additional memory to store indices.

LO-regularization

Directly solving for LO-regularization is difficult as it should try all possible subsets of variables.



Are there more efficient ways to obtain sparse models?

L1 regularization

Convex relaxation of L0 regularization.



As the dimensionality of w increases, the norm ball will have increasingly more number of corners.

Results in parameter selection

L2- vs L1-regularization

L2 regularization promotes **grouping** – results in equal weights for correlated features L1 regularization promotes **sparsity** – selects few informative features

L2-regularization





L1-regularization

Lp-Norms

General case $\|w\|_p$ where p can be any non-negative number



Lp norms with p < 2 promotes sparsity Lp norms with p > 2 promotes grouping

Example: Linear Regression

Fit a linear model *w*, that minimizes the residual between the observed and predicted values

$$\min_{w} \frac{1}{N} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{w} \|_{2}^{2}$$
$$\boldsymbol{w} = (\boldsymbol{X}^{T} \boldsymbol{X})^{-1} \boldsymbol{X}^{T} \boldsymbol{y}$$



Ridge Regression

Linear regression with L2-regularization



Shrinks all variables to zero – reduces variance while introducing bias.

Lasso (L1-regularized Linear Regression)

Shorthand for Least Absolute Shrinkage and Selection Operator Linear regression with L1-regularization

$$\min_{w} \frac{1}{N} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{w} \|_{2}^{2} + \lambda \Omega(\boldsymbol{w})$$
$$\Omega(\boldsymbol{w}) = \| \boldsymbol{w} \|_{1}$$



Not all variables become zero at the same time – results in variable selection

L1-Regularization for General Loss

L1-regularization can be coupled with various types of loss to obtain sparse models e.g.) logistic regression, SVM

 $\min_{w} l(X, y, w) + \lambda \Omega(w)$

$$l(\mathbf{X}, \mathbf{y}, \mathbf{w}) = \frac{1}{N} \sum_{i=1}^{N} [(1 - y_i) \mathbf{w}^T \mathbf{x}_i + \log(1 + \exp(-\mathbf{w}^T \mathbf{x}_i))]$$

 $\Omega(\boldsymbol{w}) = \|\boldsymbol{w}\|_1 \qquad \|\boldsymbol{w}\|_1 \leq \lambda$

How can we then **solve** for such l1-regularized objectives?

Gradient Descent

Determine the descent direction as the gradient at the point, determine the step size t, and then move to that direction.

given a starting point $x \in \operatorname{dom} f$. repeat

1. $\Delta x := -\nabla f(x)$.

2. Line search. Choose step size t via exact or backtracking line search.

3. Update. $x := x + t\Delta x$.

until stopping criterion is satisfied.

Repeat this process until a stopping criterion is met.

Cannot be applied to optimization of non-smooth objectives



Subgradient Method

Compute a set of gradient, defining the gradient on non-smooth points



Subgradient Method

A vector g is a subdifferential of f at x, if $f(y) = f(x) + d^T(y - x), \forall y$



Suffers from slow convergence, and solutions obtained are usually non-sparse

Proximal Gradient

Specifically tailored for regularized optimization problems where g(w) is differentiable but h(w) is a general convex function.





The proximal operator, prox(w) should be efficient – closed form solution preferred

Proximal Operator for L1-norm regularization

Iterative soft-thresholding operator – reduce the absolute value of each component of *w* by lambda, and if the resulting value is below zero, set it to zero.



Coordinate Descent

Optimize one variable at a time while fixing all others.

$$\min_{w_j} \nabla_j f(w^t) (w_j - w_j^t) + \frac{1}{2} \nabla_{jj}^2 f(w^t) (w_j - w_j^t) + \lambda |w_j|$$
$$w_j^* = prox_{\frac{\lambda}{\nabla_{jj}f}|\cdot|} \left(w_j^t - \frac{\nabla_j f(w_j^t)}{\nabla_{jj}^2 f} \right)$$

w_j^* can be obtained by solving for the loss with coordinate j and soft-thresholding the solution.

[Bach et. al] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, Optimization with Sparsity Inducing Penalties, F & T in Machine Learning, 2012

Stochastic Subgradient Descent

Subgradient method using noisy unbiased subgradients for scalable optimization

step size at iteration t $w_{t+1} = w_t - a_t \widetilde{g_t}$ gradient direction at iteration t

$$E(\widetilde{g_t}|\boldsymbol{w}_t) = g_t \in \partial f(\boldsymbol{w}_t)$$



A random vector is a noisy unbiased subgradient for $f: R \to R$ at x, if for all z $f(z) \ge f(w) + (E\tilde{g})^T (z - w)$

Define the optimal value as $f^* = \min\{f(w_1), \dots, f(w_t)\}$

Scales and works well for ML problems

However, SGD often does not generate sparse solutions for l1-regularized objectives

Regularized Dual Averaging Method

Consider all past subgradients of the loss and the whole regularization term when solving for a regularized objective

$$w_{t+1} = \underset{w}{\operatorname{argmin}} \left\{ \langle \overline{g_t}, w \rangle + \Omega(w) + \frac{\beta_t}{t} h(w) \right\}^{\operatorname{auxiliary Strongly}}_{\operatorname{convex function}}$$

$$\overline{g_t} = \frac{t-1}{t} \overline{g_{t-1}} + \frac{1}{t} g_t$$
Average of all gradients over t iterations
$$|g_t^i| \le \lambda$$

$$w_{t+1}^i = \begin{cases} 0, & |g_t^i| \le \lambda \\ -\frac{\sqrt{t}}{\gamma} (g_t^i - \lambda sign(g_t^i)), & otherwise \end{cases}$$

For l1-regularization RDA has a closed form solution

[Xiao et. al] L. Xiao, Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization, JMLR 2010

Regularized Dual Averaging Method

RDA obtains sparser solutions compared to those obtained by SGD and proximal gradient



RDA has a convergence rate of $1/\sqrt{T}$ for general convex regularizers.

[Xiao et. al] L. Xiao, Dual Averaging Methods for Regularized Stochastic Learning and Online Optimization, JMLR 2010

Sparse Coding

Find sparse decomposition of each data instance x_i , using some known dictionary atoms



Sparse Coding and Dictionary Learning

Find both a sparse decomposition α_i of each data instance x_i , and the bases D



Can use alternating optimization to solve for each variable at a time.

Alternative Sparse Methods

Use greedy algorithms to directly solve for IO-norm

$$\begin{split} \min_{w} \frac{1}{N} \| \boldsymbol{y} - \boldsymbol{X} \boldsymbol{w} \|_{2}^{2} & s.t. \| \boldsymbol{w} \|_{0} \leq \lambda \\ \min_{\alpha} \| \boldsymbol{x} - \boldsymbol{D} \boldsymbol{\alpha} \|_{2}^{2} & s.t. \| \boldsymbol{\alpha} \|_{0} \leq \lambda \end{split}$$

Enforce sparsity by setting the desired number of nonzero variables.

- Matching Pursuit, Orthogonal Matching Pursuit

Matching Pursuit

At each step, select the column that is most correlated with the input

Orthogonal Matching Pursuit

At each step, select the column that helps reduce the objective the most

$$\begin{split} \min_{\alpha} \| \boldsymbol{x} - \boldsymbol{D}\boldsymbol{\alpha} \|_{2}^{2} & \text{s.t.} \| \boldsymbol{\alpha} \|_{0} \leq \lambda \\ \Gamma &= \emptyset \\ \text{for } i = 1, \dots, \lambda \text{ do} \\ i &= \arg \min_{i \in \Gamma^{c}} \left\{ \min_{\alpha'} \| \boldsymbol{x} - \boldsymbol{D}_{\Gamma \cup \{i\}} \boldsymbol{\alpha'} \|_{2}^{2} \right\} \\ \Gamma &\leftarrow \Gamma \cup \{i\} \\ r \leftarrow \left(I - \boldsymbol{D}_{\Gamma} \left(\boldsymbol{D}_{\Gamma}^{T} \boldsymbol{D}_{\Gamma} \right)^{-1} \boldsymbol{D}_{\Gamma}^{T} \right) \boldsymbol{x} \\ \alpha_{\Gamma} \leftarrow \left(\boldsymbol{D}_{\Gamma}^{T} \boldsymbol{D}_{\Gamma} \right)^{-1} \boldsymbol{D}_{\Gamma}^{T} \boldsymbol{x} \end{split}$$

All coefficients extracted so far are updated at each step

Image Restoration

Eliminate noise from the original image



Original



Restored
Image Restoration with Sparse Coding

Reconstruct the image with sparse combination of learned bases



Image Inpaiting

Works with even larger level of noise





Original

Restored

Self-taught Learning

Transfer learning when source domain is **unlabeled** and only remotely related to the target task – based on the intuition that human learning is largely unsupervised.

SourceTarget $\{x_u^{(i)}\}_{i=1}^k$ $x_u^{(i)} \in R^n, k >> m$ $\{(x_l^{(i)}, y^{(i)})\}_{i=1}^m$ $x_l^{(i)} \in R^n, y^{(i)} \in \{1, ..., T\}$





Learn a higher-level, more abstract feature for the given feature type (modality)

[Raina et al.] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-taught Learning: Transfer Learning from Unlabeled Data, ICML 2007

Given unlabeled data x, find good bases b using sparse coding and dictionary learning.



[Raina07] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-taught Learning: Transfer Learning from Unlabeled Data, ICML 2007

Use alternating optimization - Solve for each variable while fixing the other, and alternate the process until convergence (Efficient algorithm introduced in 07)



[Raina07] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-taught Learning: Transfer Learning from Unlabeled Data, ICML 2007 [Lee07] H. Lee, A. Battle, R. Raina, A. Y. Ng, Efficient Sparse Coding Algorithms, NIPS 2007

Then reconstruct the examples in the target dataset using the learned bases. Reconstruction

 $\min_{a,b} \sum_{i} \underbrace{\left| |x_{u}^{(i)} - \sum_{j} a_{j}^{(i)} b_{j} | \right|_{2}^{2}}_{\text{Bases}^{2}} + \beta \sum_{i} ||a^{(i)}||_{1} \quad ||b_{j}||_{2}^{2} \le 1, \forall j \in 1, \dots, s$ $\sum_{i} = 0.8 * \left| b_{i} \right|_{2} = 0.3 * \left| b_{i} \right|_{2} + 0.3 * \left| b_{i} \right|_{2} = 0.5 * \left$

Then train a classifier (such as a SVM) over the obtained sparse codes

[Raina et al.] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-taught Learning: Transfer Learning from Unlabeled Data, ICML 2007

Results show that information gained from unlabeled data is actually useful.

Image Classification

Handwritten Character Recognition



Method	Accuracy		
Baseline	16%		
PCA	37%		
Sparse coding	47%		



Method	Accuracy		
Raw	54.8%		
PCA	54.8%		
Sparse coding	58.5%		

[Raina et al.] R. Raina, A. Battle, H. Lee, B. Packer, A. Y. Ng, Self-taught Learning: Transfer Learning from Unlabeled Data, ICML 2007

Learning Task Grouping in Multitask Learning

Allow the learners to selectively share the information across the tasks.



Can discover true support without having to provide the number of groups

[Kumar13] A. Kumar, H. Daume III, Learning Task Grouping and Overlap in Multi-Task Learning

Sparse Feature Encoding

Soft-Encoding – encodes each feature as a sparse combination of visual words.



[Wang et. al] J. Wang, J. Yang, K. Yu, F. Lv, T. Huang, and Y. Gong, Locality-Constrained Linear coding for Image Classification, CVPR 2010

Learning Sparse Filters for a CNN

Use sparse coding to reduce the number of convolutional kernels



[Liu et. al] B. Liu, M. Wang, H. Foroosh, M. Tappen, and M. Penksy, Sparse Convolutional Neural Networks, CVPR 2015

Learning Sparse Weights for a CNN

Learn sparse network weights to reduce memory usage



	ImageNet			
Model	Top-1	Top-5	Memory	
Reported [11]	59.3%	81.8%	-	
Caffe Version [9]	57.28%	80.44%	233 MB	
Sparse (ours)	55.60%	80.40%	58 MB	

Summary

Sparsity regularizations result in learning a model with *few nonzero* parameters. It is beneficial for *feature selection, model analysis*, and *model compression*.

L1-regularization uses 1-norm for regularization, which enduces sparsity. Since 1norm is *non-differentiabl*e, we can solve it through optimization methods such as *subgradient descent, proximal gradient*, and *regularized dual averaging method*.

Sparse coding encodes an input variable as a *sparse combination of some bases*. It is useful for multiple applications such as *image restoration* and *transfer learning*.

Part 2: Structured Sparsity

- Introduction to Structured Sparsity
- Group-structured sparsity
- Block-structured sparsity
- Optimizing for (2,1)-norm
- Tree-structured sparsity
- Graph-structured sparsity
- Example: Image inpainting with hierarchical dictionary learning
- Example: Multi-task learning
- Example: Learning decorrelated attributes

Structured Sparsity

Sparsity regularizations that prefer *certain structure* among the variables over others.



Enables to exploit known structure – e.g.) To reconstruct handwritten characters, we can exploit the fact that they form connected components

(2,1)-norm

Mixed norm that has 1-norm over 2-norm groups.







Used to promote sparsity at group level

[Yuan et al.] M. Yuan and Y. Lin, Model Selection and Estimation in Regression with Grouped Variables, Journal of Royal Statistics Society, 2006

Group Sparsity

Structured sparsity with groups of variables



Group Sparsity

Use L2/L1-regularization to select / drop variables in a group

$$\min_{\boldsymbol{w}} l(\boldsymbol{X}, \boldsymbol{y}, \boldsymbol{w}) + \lambda \|\boldsymbol{w}\|_{2,1}$$

$$\|\boldsymbol{w}\|_{2,1} = \sum_{l}^{L} \|\boldsymbol{w}_{l}\|_{2}$$



Used to promote sparsity at group level

Group Lasso

Correlated variables gets selected / dropped at the same time

$$\min_{\mathbf{w}} \frac{1}{N} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_{2}^{2} + \lambda \sum_{l=1}^{L} \|\mathbf{w}_{l}\|_{2}$$



Sparse Group Lasso

Group lasso does not yield sparsity within a group. All variables selected will have non-zero values.

$$\min_{w} \frac{1}{N} \| \mathbf{y} - \mathbf{X} \mathbf{w} \|_{2}^{2} + \lambda_{1} \sum_{l=1}^{L} \| \mathbf{w}_{l} \|_{2} + \lambda_{2} \| \mathbf{w} \|_{1}$$
group lasso
sparse group lasso

Use additional I1-regularization to yield sparsity at individual feature level.

Block Sparsity

Structured Sparsity on a 2D Grid



Block Sparsity

Same as group sparsity applied to a sequence. Group each row or columns using 2norm, and apply 1-norm on the groups



Exclusive Lasso

Promote competition for the features between tasks

```
\min_{\boldsymbol{w}} l(\boldsymbol{X}, \boldsymbol{y}, \boldsymbol{w}) + \lambda \Omega(\boldsymbol{w})
```









[Zhou10] Y. Zhou, R. Jin, and S. C. H. Hoi, Exclusive Lasso for Multi-task Feature Selection, AISTAT 2010

Optimization for L2/L1-Regularized Objectives

(2,1)-norm is nonsmooth.

proximal operator

$$prox(\mathbf{w})_l = \left[1 - \frac{\lambda}{\|\mathbf{w}_l\|}\right]_+ \mathbf{w}_l$$

Regularized dual averaging

$$w_{t+1}^{i} = \begin{cases} 0, & \text{if } |g_{t}^{i}| \leq \lambda \\ -\frac{1}{\sigma} (g_{t}^{i} - \lambda sign(g_{t}^{i})), & \text{otherwise} \end{cases}$$

Can use proximal gradient and regularized dual averaging method

Block Coordinate Descent

Optimize a group of variable at a time, while fixing all others.

$$\boldsymbol{w}_{g}^{*} = prox_{\frac{\lambda}{L_{t}}} \left(\boldsymbol{w}_{g} - \frac{1}{L_{t}} \nabla_{g} f(\boldsymbol{w}_{g}^{t}) \right)$$

$$prox(\boldsymbol{w}) = \left[1 - \frac{\lambda}{\|\boldsymbol{w}\|}\right]_{+} \boldsymbol{w}$$

The solution w_{g}^{*} can be obtained through group-soft thresholding

[Blondel et al.] M. Blondel, K. Seki, K. Uehara, Block Coordinate Descent Algorithms for Large-Scale Sparse Multiclass Classification

Tree Sparsity

Structural sparsity on a tree



If a node is removed, then all its descendant nodes are dropped.

[Jenatton11a] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, Proximal Methods for Sparse Hierarchical Dictionary Learning, ICML 2010

Hierarchical Sparsity-Inducing Norm

Perform group lasso where each group G_j contains node j and all its descendants.



$$\begin{split} \Omega(\beta) &= \sum_{g \in \mathcal{G}} w_g || \boldsymbol{\beta}_g ||_2 \\ \boldsymbol{\beta}_1 &= \beta_1 \quad \boldsymbol{\beta}_2 = \beta_2 \quad \boldsymbol{\beta}_3 = \beta_3 \\ \boldsymbol{\beta}_4 &= \{\beta_1, \beta_4\} \quad \boldsymbol{\beta}_5 = \{\beta_2, \beta_3, \beta_5\} \\ \boldsymbol{\beta}_6 &= \{\beta_1, \beta_2, \beta_3, \beta_4, \beta_5\} \end{split}$$

[Jenatton11a] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, Proximal Methods for Sparse Hierarchical Dictionary Learning, ICML 2010

Tree-guided Group Lasso

Exploit a given tree structure on the output values to recover the true structure in the parameters



Idea: use overlapping groups in group lasso

[Kim et al.] S. Kim and E. P. Xing, Tree Guided Group Lasso for Multi-Task Regression with Structured Sparsity, ICML 2010

Tree-guided Group Lasso

Promotes sharing and competition between tasks in hierarchical manner.



Simplest case where there are only two outputs

[Kim et al.] S. Kim and E. P. Xing, Tree Guided Group Lasso for Multi-Task Regression with Structured Sparsity, ICML 2010

Tree-guided Group Lasso

Promotes sharing and competition between tasks in hierarchical manner.



[Kim et al.] S. Kim and E. P. Xing, Tree Guided Group Lasso for Multi-Task Regression with Structured Sparsity, ICML 2010

Graph-Structured Sparsity

Structured sparsity on a generic graph





[Hegde15] C. Hegde, P. Indyk, and L. Schmidt, A Nearly-Linear Time Framework for Graph-Structured Sparsity, ICML 2015

Graph-Guided Fused Lasso

Use known graph structure among the output to constrain correlated variables to have similar parameters



Multitask Feature Learning

Assume that there exist some latent shared features that are shared across multiple tasks, and learn them.



[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Multitask Feature Learning

Given N input features x and label $\mathcal{Y}_t\;$ for each task t, Simultaneously learn the transformation U and model parameter θ_t for each class t



1. Learn classifiers on the transformed features in shared feature space

2. Promote a common sparsity pattern in the new parameters

[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Sharing features via Sparsity Regularization

Using group sparsity regularization, enforce each learner to use features that are informative across multiple tasks.



[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Convex Multitask Feature Learning

Previous formulation is non-smooth, which makes it challenging to solve. Thus we solve an equivalent form instead. \rightarrow Replace features with a covariance matrix Ω that measures the relative effectiveness of each dimension

 w_t : model parameter for task t, Ω : covariance matrix on original features.

$$W^*, \Omega^* = \arg \min \sum_t \sum_n \ell(w_t^T x_n, y_{nt})$$
 Model prediction loss in
the **original** feature space
 $+ \gamma \sum_t w_t^T \Omega^{-1} w_t + \gamma \epsilon \operatorname{Trace}(\Omega^{-1})$
Trace norm
(sum of the diagonal entries for a PSD matrix)

Multitask Feature Learning - Result

Dataset

Animals with Attributes (AWA) 30,475 images, 50 classes, 85 attributes



Outdoor Scene Recognition (OSR) 2,688 images, 8 classes, 6 attributes


Multitask Feature Learning - Result

No Sharing (Visual Features) > Attributes-based Prediction



No sharing-Obj.: Independent SVMs trained on visual features No sharing-Attr. : Object recognition on predicted attributes as in [Lampert09].

[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Multitask Feature Learning - Result

Shared Features (Object & Attributes) > Shared Features (Objects) > No Sharing



Sharing-Obj.: Multitask Feature Learning on object classifiers Sharing-Attr. : Multitask Feature Learning on object + attribute classifiers

[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Multitask Feature Learning - Result

Predicted object categories Red: incorrect prediction

Grizzly Bear

More robustness to background clutter from features refined with attributes





Ours

Ours



No Sharing Polar Bear

Ours

Dalmatian

Even when our method fails, it often makes more semantically "close" predictions.





Rhinoceros



Wolf

Cow

[Hwang11a] S. J. Hwang, F. Sha, K. Grauman, Sharing Features between Objects and Their Attributes, CVPR 2011

Image Inpainting with Hierarchical Sparse Coding

Use hierarchical sparsity inducing-norm to reconstruct an input image with a hierarchy of image patches



Hierarchical Sparse Coding

Use proximal operator to solve for the hierarchical sparsity-inducing norm

$$\begin{split} \min_{\mathbf{D}\in\mathcal{D},\mathbf{A}\in\mathcal{A}} \frac{1}{n} \sum_{i=1}^{n} \left[\frac{1}{2} ||\mathbf{x}^{i} - \mathbf{D}\boldsymbol{\alpha}^{i}||_{2}^{2} + \lambda \Omega(\boldsymbol{\alpha}^{i}) \right] \\ \Omega(\boldsymbol{\alpha}) &\triangleq \sum_{g\in\mathcal{G}} w_{g} ||\boldsymbol{\alpha}_{|g}|| \\ \min_{\mathbf{v}\in\mathbb{R}^{p}} \frac{1}{2} ||\mathbf{u} - \mathbf{v}||_{2}^{2} + \lambda \Omega(\mathbf{v}) \end{split}$$

Problem – No closed form solution exists for group lasso with overlapping groups

[Jenatton11a] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, Proximal Methods for Sparse Hierarchical Dictionary Learning, ICML 2010

Hierarchical Sparse Coding

Use primal-dual approach to solve for the dual problem

Primal
$$\begin{aligned} \min_{\mathbf{v}\in\mathbb{R}^{p}} \frac{1}{2} \|\mathbf{u} - \mathbf{v}\|_{2}^{2} + \lambda\Omega(\mathbf{v}) \\ \max_{\boldsymbol{\xi}\in\mathbb{R}^{p\times|\mathcal{G}|}} -\frac{1}{2} \left(\|\mathbf{u} - \sum_{g\in\mathcal{G}} \boldsymbol{\xi}^{g}\|_{2}^{2} - \|\mathbf{u}\|_{2}^{2}\right) \\ \text{S.t. } \forall g\in\mathcal{G}, \ \|\boldsymbol{\xi}^{g}\|_{*} \leq \lambda w_{g} \text{ and } \boldsymbol{\xi}_{j}^{g} = 0 \text{ if } j \notin g \end{aligned}$$

Algorithm 1 Block coordinate ascent in the dualInputs: $\mathbf{u} \in \mathbb{R}^p$ and set of groups \mathcal{G} .Outputs: $(\mathbf{v}, \boldsymbol{\xi})$ (primal-dual solutions).Initialization: $\mathbf{v} = \mathbf{u}, \boldsymbol{\xi} = 0$.while (maximum number of iterations not reached) dofor $g \in \mathcal{G}$ do $\mathbf{v} \leftarrow \mathbf{u} - \sum_{h \neq g} \boldsymbol{\xi}^h$. $\boldsymbol{\xi}^g \leftarrow \Pi^*_{\lambda w_g}(\mathbf{v}_{\lg})$.end forend while $\mathbf{v} \leftarrow \mathbf{u} - \sum_{g \in \mathcal{G}} \boldsymbol{\xi}^g$.

Convex constraints for each vector ξ^{g} are separable

- Can be efficiently solved using block-coordinate ascent

[Jenatton11a] R. Jenatton, J. Mairal, G. Obozinski, and F. Bach, Proximal Methods for Sparse Hierarchical Dictionary Learning, ICML 2010

Image Inpainting Results

Hierarchical dictionary learning yields much less reconstruction noise compared to flat dictionary learning with flat l1-regularization.



L1



Free- spar	sity
-------------------	------

Noise	50%	90%
Flat	19.3	72.1
Hierarchical	18.6	65.9

Problem: difficult to distinguish between co-occurring attributes



Forest animal? Brown? Has ears? Combinations?

Motivation: learning correct attributes is crucial for applications such as image search and zero-shot learning.



Solution: promote competition between



Attributes can be grouped into multiple semantic categories

Texture	Character	Color	Parts	Activity	Nutrition	Habitat
patches spots stripes furry hairless toughskin	fierce timid smart group solitary nestspot domestic	black white blue brown gray orange red	flippers chewteeth hands meatteeth hooves buckteeth pads strainteeth paws horns longleg claws longneck tusks tail bipedal	flys hops swims tunnels walks fast slow	fish meat plankton vegetation insects forager grazer	coastal desert bush plains forest fields jungle
Behavior		yellow	quadrapedal Shape	weak	hunter scavenger skimmer	mountains ocean
active inactive nocturnal	hibernate agility	b sn	ig bulbous nall lean	musere	stalker	water tree cave

Promote sharing between attributes in the same group, while promoting competition for features in different groups



Decorrelation model predicts attributes much better in unusual cases.



Not brown No eye Not boxy No mouth No ear underparts

Brown wing Olive back Crested head

This model can also accurately localize the part-based attributes.

Standard Ours

[Jayaraman et al.] D. Jayaraman, F. Sha, and K. Grauman, Decorrelating Visual Attributes by Resisting Urge to Share, CVPR 2014

Blue back

SPAMS (SPArse Modeling Software)

A popular optimization toolbox for solving most of the methods introduced in this presentation

http://spams-devel.gforge.inria.fr/

- Implements OMP, Lasso, LARS, coordinate descent and proximal methods for I1regularization.
- Provides proximal toolbox for solving I2/I1-reg., sparse group lasso, treestructured regularization, and structured sparsity with overlapping groups.

[Mairal et al.] J. Mairal, F. Bach, and J. Ponce, Sparse Modeling for Image and Vision Processing, Foundations and Trends in Computer Graphics and Vision, 2014

Summary

Structured sparsity is a kind of sparsity that prefers *certain structure* among the variables, based on some *prior knowledge*.

There are various types of structured sparsity, including *group sparsity*, *block sparsity*, *tree-structured sparsity*, and *graph-based sparsity*. Structured sparsity is often enforced as *mixed norm* regularization, such as (2,1)-norm

Structured sparsity is useful for multiple applications, including *multi-task learning*, and hierarchical dictionary learning for *image denoising*.