

Introduction to Sparse Filtering

May 27, 2015

Aim of this presentation

In this presentation we are going to discuss about *sparse filtering*.

Sparse filtering is a specific algorithm for unsupervised learning, but it provides the opportunity to discuss aspects of unsupervised learning that are general and conceptually stimulating.

Intuition on Sparse Filtering

SF is a simple *algorithm* or *learning module* to perform *unsupervised feature distribution learning* proposed by Ngiam et al. in 2011 [5].

- Finality or domain of applicability: *unsupervised learning*;
- Modality or approach: *feature distribution learning*;
- Desiderata: *simplicity* (i.e.: few hyperparameters) and *efficiency* (i.e.: scalable wrt the dimensionality of the input).

Unsupervised Learning (I)

Unsupervised Learning means learning from data \mathcal{D} without external information.

Learn a good model generating *representation* of data

$$f : \mathcal{D} \rightarrow \mathcal{R}$$

Unsupervised learning is a underdetermined problem.

Unsupervised Learning (II)

What do we learn in absence of external guidance?

Two main approaches [5]:

- *Data Distribution Learning*: learn the *true* distribution of the data \mathcal{D} .
- *Feature Distribution Learning*: learn a *useful* distribution of the representations \mathcal{R} .

Data Distribution Learning

Data Distribution Learning is the traditional approach to unsupervised learning in which, given data \mathcal{D} , we try to model the distribution of the process that generated \mathcal{D} .

Several mainstream algorithms: *Boltzmann machines*, *autoencoders*, *independent component analysis* [5].

Implicit assumption: learning the *true structure of the data* (i.e.: the statistical description of the process generating the data) will automatically provide a *useful* representation.

Feature Distribution Learning (I)

Feature Distribution Learning is an innovative approach to unsupervised learning in which, given data \mathcal{D} , we try to model the distribution of the representation \mathcal{R} in order to maximize its usefulness.

SF being the first algorithm of this kind [5].

Implicit assumption: some forms of representation are better than others and they will automatically provide a *useful* representation.

Feature Distribution Learning (II)

Assuming the conceptual framework of feature distribution learning we may now wonder:

- ① *What sort of feature distribution may we want to learn?*
- ② *How do we learn a feature distribution?*
- ③ *Is feature distribution learning feasible at all?*

Sparsity

1. *What sort of feature distribution may we want to learn?*

A *sparse* distribution, that is a distribution where most of the values are zero.

- *Practical reason:* sparse representation proved successful in many machine learning task (e.g.: *sparse deep belief networks* [7] or *k-sparse autoencoders* [4]);
- *Analogical reason:* biological systems implements sparse distributed representations (e.g.: modelling V1 cortex coding [6]);
- *Formal reason:* sparse distribution has low entropy¹ ([1])

¹how important is this?

Sparse Filtering

2. *How do we learn a feature distribution?*

SF is an *algorithm* or *learning module* to perform *unsupervised feature distribution learning* that generates *sparse representations*.

Sparse Filtering

SF is an *algorithm or learning module* to perform *unsupervised feature distribution learning* that generates *sparse representations*.

Given a dataset²:

$$\text{raw features} \left\{ \underbrace{\begin{bmatrix} .3 & .4 & .3 & \dots & .7 \\ .2 & .7 & .3 & \dots & .3 \\ .3 & .8 & .5 & \dots & .6 \\ \dots & \dots & \dots & \dots & \dots \\ .2 & .1 & .8 & \dots & .4 \end{bmatrix}}_{\text{samples}} \right\} \xrightarrow{SF} \left\{ \underbrace{\begin{bmatrix} 0 & 0 & 0 & \dots & .7 \\ 0 & 0 & 0 & \dots & .6 \\ 0 & .7 & 0 & \dots & 0 \\ 0 & .8 & 0 & \dots & 0 \\ .9 & 0 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & .8 & \dots & 0 \end{bmatrix}}_{\text{samples}} \right\} \text{SF features}$$

²notice the slightly unusual convention of having features along the rows and samples along the columns

SF - Sparsity

We achieve sparsity enforcing three properties:

- 1 *Population Sparsity*: each sample has few non-zero values;
- 2 *Lifetime Sparsity*: each feature has few non-zero values;
- 3 *High Dispersal*: activity on each row should be constant.

$$\begin{bmatrix} 0 & 0 & 0 & \dots & .7 \\ .7 & 0 & 0 & \dots & .6 \\ 0 & .8 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & .8 & \dots & 0 \end{bmatrix}$$

SF Algorithm

Minimize the following *loss function*

$$\operatorname{argmin}_W \left\| \left\| \left\| f(WX) \right\|_{L2, \text{row}} \right\|_{L2, \text{column}} \right\|_{L1}$$

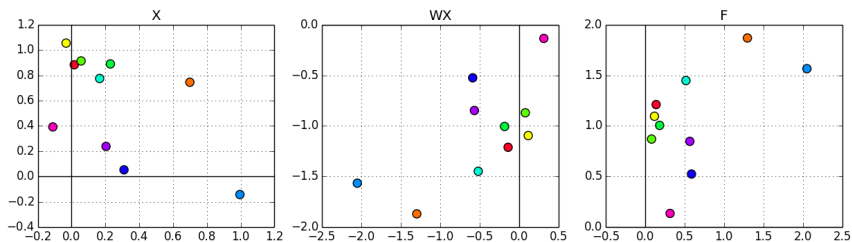
through *gradient descent*.

This ugly formula can be decomposed into four intuitive steps.

SF Algorithm - Step 1

Non-linear processing:

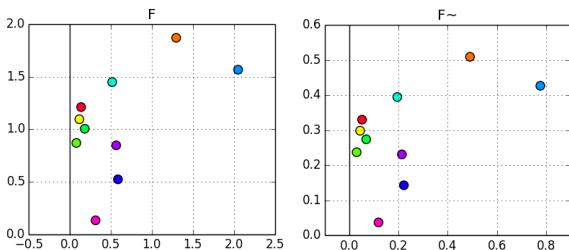
$$F = f(WX) = |WX|$$



SF Algorithm - Step 2

Normalization along the rows (features):

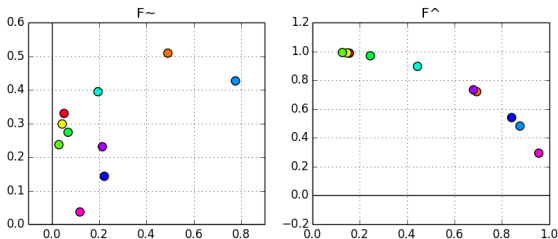
$$\tilde{F} = \frac{F}{\|F\|_{L2, row}}$$



SF Algorithm - Step 3

Normalization along the columns (samples):

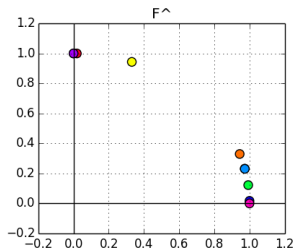
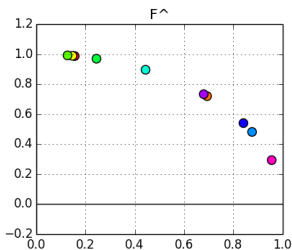
$$\hat{F} = \frac{\tilde{F}}{\|\tilde{F}\|_{L2, column}}$$



SF Algorithm - Step 4

Minimization of L1 norm:

$$\|\hat{F}\|_{L1}$$



SF Algorithm - Observations

- 1 *Population sparsity* is achieved by minimizing L1 norm;
 - 2 *High dispersal* is enforced by row normalization by imposing the same mean square activation for each feature;
 - 3 *Lifetime sparsity*: follows from the previous properties.
- These properties are imposed on the *feature distribution* not on the *data distribution*.
 - Is the intuition that we are performing a sort of *constrained scattering* correct?

Original Results

3. *Is feature distribution learning feasible at all?*

Results in [5]:

- *Timing and scaling*: SF shown to scale better than ICA, sparse coding or sparse autoencoders;
- *Processing of natural images*: SF learns Gabor-like filters;
- *Object classification on STL-10*: SF representations allows a linear SVM to achieve better performances than raw, random weights, k-means, and ICA representations;
- *Phoneme classification on TIMIT*: SF representations allows a linear or RBF SVM to achieve state-of-the-art performances.

Evaluating SF: Pros

- ✓ State-of-the-art performances
- ✓ Neat mathematical formulation
- ✓ Hyperparameter-light
- ✓ Highly computationally scalable
- ✓ Stackable
- ✓ Extensible³

³Can we use SF as a starting point to develop different feature distribution learning algorithms?

Evaluating SF: Cons

- × Data structure-agnostic
- × Fragility⁴
- × Sensitivity to initialization

⁴Could the extension of SF compromise the aim of learning a sparse distribution?

Discussion (I)

- *Is soft-absolute the best non-linearity?*
[5] suggests the potential use of other functions, but no studies are available.
- *Is L2 normalization the best normalization?*
It may be possible to project on other surfaces than a hypersphere, but no studies are available.
- *Can we earn something from stacking SF?*
[5] run tests to study the filter-like behaviour of stacked SF, but no further results are reported.

Discussion (II)

- *When does the algorithm fail?*
Testing shows that sometimes SF just fail in finding a sparse distribution.
- *How does failure relate to initialization?*
Testing shows that radically different solutions are achieved with radically different initialization.
- *Can we prevent failure?*
SF may be improved if the likelihood of failure could be decreased, or if failure-bound instances may be stopped earlier.
- *What is the optimal stopping criterion?*
Literature show that the performances of SF strongly depends on the number of iterations.

Follow-ups: Practical Application

Thaler (1/218) [2] and *Romaszko* (6/218) [8] used SF successfully in the *Kaggle Black Box Challenge*.

Performance-oriented application of the SF algorithm in a classification pipeline.

Follow-ups: Practical Application

*Ngiam et al.*⁵

Research approach

Shallow architecture (1 layer)

Overcomplete representation

Processing all the data in an
unsupervised scheme

*Thaler*⁶

Applicative approach

Deep architecture (2 layer)

Undercomplete representation

Training on training data and
processing on testing data

⁵Some of these details are inferred from [5]

⁶These details are based on the technical report released by Thaler

Follow-ups: Comparisons and Extensions

- *Yang et al.* [10] extends the SF algorithm adding L2 weight regularization;
- *Zhang et al.* [11] compares different algorithms for unsupervised learning (including SF);
- *Lederer et al.* [3] suggests a connection between SF and *random matrix theory*;
- *Romero et al.* [9] improves on learning sparsity.

Some Further Questions

- *How to improve the sparsity learning algorithm?*
Finding optimal ways to deal with failure and weight initialization.
- *How to apply sparse filtering in a train&test scenario?*
Representation produced by SF depends on the other samples in the batch we are normalizing.
- *Can we use SF in a semi-supervised scenario?*
SF may be used to learn representation out of big unlabelled datasets, before performing supervised or weakly supervised learning.
- *Can we use SF in a semi-supervised scenario where we aim at learning disentanglement?*

(Emotional) Disentangling Sparse Filtering

$$F = f(WX) = \begin{bmatrix} A & B \\ C & D \end{bmatrix} \text{ where } \underbrace{\begin{bmatrix} A \\ C \end{bmatrix}}_{\text{emo sample}} \quad \underbrace{\begin{bmatrix} B \\ D \end{bmatrix}}_{\text{nemo samples}} \begin{matrix} \text{nemo features} \\ \text{emo features} \end{matrix}$$

$$\mathcal{L} = \left\| \left\| \left\| \begin{bmatrix} A & B \\ C & D \end{bmatrix} \right\|_{L2, \text{row}} \right\|_{L2, \text{column}} \right\|_{L1} + \lambda_D \|D\|_{L1}$$

$$\mathcal{L} = \left\| \left\| \left\| \begin{bmatrix} A & B \\ C & D \end{bmatrix} \right\|_{L2, \text{row}} \right\|_{L2, \text{column}} \right\|_{L1} + \lambda_D \|D\|_{L1} + \lambda_A \|A\|_{L1}$$

Sparse Filtering Off-the-shelf

- *Matlab Implementation:*
<https://github.com/jngiam/sparseFiltering>
- *Python Porting:*
<https://github.com/martinblom/py-sparse-filtering>

References I

- [1] Chakra Chennubhotla and Allan Jepson.
Sparse coding in practice.
In Proc. of the Second Int. Workshop on Statistical and Computational Theories of Vision, 2001.
- [2] Ian J. Goodfellow, Dumitru Erhan, Pierre Luc Carrier, Aaron Courville, Mehdi Mirza, Ben Hamner, Will Cukierski, Yichuan Tang, David Thaler, Dong-Hyun Lee, Yingbo Zhou, Chetan Ramaiah, Fangxiang Feng, Ruifan Li, Xiaojie Wang, Dimitris Athanasakis, John Shawe-Taylor, Maxim Milakov, John Park, Radu Ionescu, Marius Popescu, Cristian Grozea, James Bergstra, Jingjing Xie, Lukasz Romaszko, Bing Xu, Zhang Chuang, and Yoshua Bengio.

References II

Challenges in representation learning: A report on three machine learning contests.

In Proceedings of the International Conference on Neural Information Processing, 2013.

[3] Johannes Lederer and Sergio Guadarrama.

Compute less to get more: Using orc to improve sparse filtering.

arXiv preprint arXiv:1409.4689, 2014.

[4] Alireza Makhzani and Brendan Frey.

k-sparse autoencoders.

arXiv preprint arXiv:1312.5663, 2013.

References III

- [5] Jiquan Ngiam, Zhenghao Chen, Sonia A Bhaskar, Pang W Koh, and Andrew Y Ng.
Sparse filtering.
In Advances in Neural Information Processing Systems, pages 1125–1133, 2011.
- [6] Bruno A Olshausen and David J Field.
Sparse coding with an overcomplete basis set: A strategy employed by v1?
Vision research, 37(23):3311–3325, 1997.
- [7] Marc'Aurelio Ranzato, Y-Lan Boureau, and Yann LeCun.
Sparse feature learning for deep belief networks.
In Proceedings of Neural Information Processing Systems, 2007.

References IV

- [8] Lukasz Romaszko.
A deep learning approach with an ensemble-based neural network classifier for black box icml 2013 contest.
In Workshop on Challenges in Representation Learning, ICML, 2013.
- [9] Adriana Romero, Petia Radeva, and Carlo Gatta.
No more meta-parameter tuning in unsupervised sparse feature learning.
arXiv preprint arXiv:1402.5766, 2014.

References V

- [10] Zhao Yang, Lianwen Jin, Dapeng Tao, Shuye Zhang, and Xin Zhang.

Single-layer unsupervised feature learning with l_2 regularized sparse filtering.

In *Signal and Information Processing (ChinaSIP), 2014 IEEE China Summit & International Conference on*, pages 475–479. IEEE, 2014.

- [11] Shaohua Zhang, Hua Yang, and Zhouping Yin.

Performance evaluation of typical unsupervised feature learning algorithms for visual object recognition.

In *Intelligent Control and Automation (WCICA), 2014 11th World Congress on*, pages 5191–5196. IEEE, 2014.