Exploiting Statistical Dependencies in Sparse Representations*



Michael Elad

The Computer Science Department The Technion – Israel Institute of technology Haifa 32000, Israel

* Joint work with





Tomer Faktor & Yonina Eldar



January 6-7, 2011 London, UK

Sparse Models, Algorithms and Learning for Large-scale data



This research was supported by the European Community's FP7-FET program SMALL under grant agreement no. 225913





Part Motivation for Modeling Dependencies



Sparse & Redundant Representations

- □ We model a signal <u>x</u> as a multiplication of a (given) dictionary **D** by a sparse vector $\underline{\alpha}$.
- □ The classical (Bayesian) approach for modeling $\underline{\alpha}$ is as a Bernoulli-Gaussian RV, implying that the entries of $\underline{\alpha}$ are iid.
- \Box Typical recovery: find $\underline{\alpha}$ from

 $\underline{\mathbf{y}} = \underline{\mathbf{x}} + \underline{\mathbf{e}} = \mathbf{D}\underline{\alpha} + \underline{\mathbf{e}}$.

This is done with pursuit algorithms (OMP, BP, ...) that assume this iid structure.

□ Is this fair? Is it correct?





Lets Apply a Simple Experiment

- □ We gather a set of 100,000 randomly chosen image-patches of size 8×8, extracted from the image 'Barbara'. These are our signals $\{\underline{y}_k\}_{\mu}$.
- □ We apply sparse-coding for these patches using the OMP algorithm (error-mode with σ =2) with the redundant DCT dictionary (256 atoms). Thus, we approximate the solution of

$$\underline{\hat{\alpha}}_{k} = \min_{\underline{\alpha}} \left\| \underline{\alpha} \right\|_{0} \text{ s.t. } \left\| \underline{y}_{k} - \mathbf{D} \underline{\alpha} \right\|_{2}^{2} \leq N\sigma^{2}$$

Lets observe the obtained representations' supports. Specifically, lets check whether the iid assumption is true.





The Experiment's Results (1)

□ The supports are denoted by <u>s</u> of length K (256), with {+1} for on-support element and {-1} for off-support ones.

$$\begin{bmatrix} \mathbf{S}_1 & \mathbf{S}_2 & \mathbf{S}_3 \\ \mathbf{S}_1 & \mathbf{S}_2 & \mathbf{S}_3 \end{bmatrix} \bullet \bullet \bullet \begin{bmatrix} \mathbf{S}_{\mathsf{K}-1} & \mathbf{S}_{\mathsf{K}} \end{bmatrix} \in \{-1, 1\}^k$$

First, we check the probability of each atom participate in the support.

As can be seen, these probabilities are far from being uniform, even though OMP gave all atoms the same weight.





The Experiment's Results (2)

Next, we check whether there are dependencies between the atoms, by evaluating

$$\frac{P\left(s_{i}=1,s_{j}=1\right)}{P\left(s_{i}=1\right)P\left(s_{j}=1\right)}, \quad 1 \leq i,j \leq K$$

This ratio is 1 if s_i and s_j are independent, and away from 1 when there is a dependency.
 After abs-log, a value away from zero indicates a dependency.



As can be seen, there are some dependencies between the elements, even though OMP did not promote them.



The Experiment's Results (3)

- We saw that there are dependencies within the support vector. Could we claim that these form a block-sparsity pattern?
- If s_i and s_j belong to the same block, we should get that

 $P\left(s_{_{i}}=1 \mid s_{_{j}}=1\right)=1$

Thus, after abs-log, any value away from 0 means a non-block structure.



We can see that the supports do not form a block-sparsity structure, but rather a more complex one.



Our Goals

- □ We would like to (statistically) model these dependencies somehow.
- \Box We will adopt a generative model for $\underline{\alpha}$ of the following form:

Generate the support vector s (of length K) by drawing it from the distribution P(s)



Generate the non-zero values for the on-support entries, drawn from the distributions $N(0, \sigma_i^2)$.

- What P(s) is? We shall follow recent work and choose the Boltzmann Machine [Garrigues & Olshausen '08]
 [Cevher, Duarte, Hedge & Baraniuk '09].
- □ Our goals in this work are:
 - $\circ~$ Develop general-purpose pursuit algorithms for this model.
 - $\circ~$ Estimate this model's parameters for a given bulk of data.



Part II Basics of the Boltzmann Machine



BM – Definition

Recall that the support vector, <u>s</u>, is a binary vector of length K, indicating the on- and the off-support elements:

$$\begin{bmatrix} \mathbf{S}_1 & \mathbf{S}_2 & \mathbf{S}_3 \end{bmatrix} \bullet \bullet \bullet \begin{bmatrix} \mathbf{S}_{K-1} & \mathbf{S}_{K} \end{bmatrix} \in \{-1, 1\}^K$$

□ The Boltzmann Machine assigns the following probability to this vector:

$$\mathsf{P}(\underline{s}) = \frac{1}{\mathsf{Z}(\mathsf{W},\underline{b})} \exp\left\{\underline{b}^{\mathsf{T}}\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}\mathbf{W}\underline{s}\right\}$$

Boltzmann Machine (BM) is a special case of the more general Markov-Random-Field (MRF) model.



The Vector <u>b</u>

 \Box If we assume that **W**=**O**, the BM becomes

$$P(\underline{s}) = \frac{1}{Z} \exp\left\{\underline{b}^{\mathsf{T}}\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}\mathbf{W}\underline{s}\right\}$$
$$= \frac{1}{Z} \exp\left\{\underline{b}^{\mathsf{T}}\underline{s}\right\} = \frac{1}{Z} \prod_{k=1}^{K} \exp\left(b_{k}s_{k}\right)$$

 \Box In this case, the entries of <u>s</u> are independent.

□ In this case, the probability of $s_k = +1$ is given by

$$p_{k} = P(s_{k} = 1) = \frac{exp(b_{k})}{exp(b_{k}) + exp(-b_{k})}$$

Thus, p_k becomes very small (promoting sparsity) if $b_k << 0$.



The Matrix W

 \Box W introduces the inter-dependencies between the entries of <u>s</u>:

 \Box Zero value W_{ij} implies an 'independence' between s_i and s_j.

- \Box A positive value W_{ij} implies an 'excitatory' interaction between s_i and s_j .
- \Box A negative value W_{ij} implies an 'inhibitory' interaction between s_i and s_j .
- \Box The larger the magnitude of W_{ij} , the stronger are these forces.
- □ W is symmetric, since the interaction between i-j is the same as the one between j-i.

$$\square \text{ We assume that} \underset{\text{diag}(\mathbf{W})=\underline{0}, \text{ since}}{\text{liag}(\mathbf{W})=\underline{0}, \text{ since}} P(\underline{s}) = \frac{1}{Z} \exp\left\{\underline{b}^{\mathsf{T}}\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}\mathbf{W}\underline{s}\right\}$$
$$= \frac{1}{Z} \exp\left\{\underline{b}^{\mathsf{T}}\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}(\mathbf{W} - \text{diag}(\mathbf{W}))\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}\text{diag}(\mathbf{W})\underline{s}\right\}$$



Special Cases of Interest

- □ When W=0 and <u>b</u>=c¹, we return to the iid case we started with. If c<<0, this model promotes sparsity.
- When W is block-diagonal (up to permutation) with very high and positive entries, this corresponds to a block-sparsity structure.
- A banded W (up to permutation) with 2L+1 main diagonals corresponds to a chordal graph that leads to a decomposable model.
- A tree-structure can also fit this model easily.





Part II MAP/MMSE Sparse Recovery Revisited



The Signal Model



D is fixed and known.

 \Box We assume that $\underline{\alpha}$ is built by:

- Choosing the support <u>s</u> with BM probability P(<u>s</u>) from all the 2^κ possibilities Ω.
- Choosing the $\underline{\alpha}_s$ coefficients using independent Gaussian entries $N(0, \sigma_i^2) \rightarrow$ we shall assume that σ_i are the same (= σ_{α}).

 \Box The ideal signal is $\underline{\mathbf{x}} = \mathbf{D}_{\underline{\alpha}} = \mathbf{D}_{\underline{s}} \underline{\alpha}_{\underline{s}}$.



The p.d.f. $P(\underline{\alpha})$ and $P(\underline{x})$ are clear and known



Adding Noise



Noise Assumed:

The noise <u>e</u> is additive white Gaussian vector with probability $P_F(\underline{e})$

$$P(\underline{y}|\underline{x}) = C \cdot exp\left\{-\frac{\left\|\underline{x} - \underline{y}\right\|^{2}}{2\sigma_{e}^{2}}\right\}$$

The conditional p.d.f.'s $P(\underline{y}|\underline{\alpha})$, $P(\underline{\alpha}|\underline{y})$, and even $P(\underline{y}|\underline{s})$, $P(\underline{s}|\underline{y})$, are all clear and well-defined (although they may appear nasty).



Pursuit via the Posterior



- □ There is a delicate problem due to the mixture of continuous and discrete PDF's. This is why we estimate the support using MAP first.
- □ MAP estimation of the representation is obtained using the oracle's formula.



The Oracle

$$\mathsf{P}\left(\underline{\alpha} \mid \underline{y}, \underline{s}\right) = \mathsf{P}\left(\underline{\alpha}_{s} \mid \underline{y}\right)$$

$$P(\underline{y} \mid \underline{\alpha}_{s}) \propto \exp\left\{-\frac{\left\|\underline{y} - \mathbf{D}_{s}\underline{\alpha}_{s}\right\|^{2}}{2\sigma_{e}^{2}}\right\}$$

 $\sigma_e^2 \sigma_a^2$

α

$$\mathsf{P}(\underline{\alpha}_{s}) \propto \mathsf{exp}\left\{-\frac{\left\|\underline{\alpha}_{s}\right\|_{2}^{2}}{2\sigma_{\alpha}^{2}}\right\}$$

$$\mathsf{P}\left(\underline{\alpha}_{s} \mid \underline{y}\right) \propto \exp\left\{-\frac{\left\|\underline{y} - \mathbf{D}_{s}\underline{\alpha}_{s}\right\|^{2}}{2\sigma_{e}^{2}} - \frac{\left\|\underline{\alpha}_{s}\right\|_{2}^{2}}{2\sigma_{\alpha}^{2}}\right\}$$

 $\mathbf{D}_{s}^{\mathsf{T}} \mathbf{y}$

- This estimate is both the MAP and MMSE.
- The oracle estimate of \underline{x} is obtained by multiplication by **D**_s.



Exploiting Statistical Dependencies in Sparse Representations By: Michael Elad

 $\underline{\hat{\alpha}}_{s} = \mathbf{D}_{s}^{\mathsf{T}}\mathbf{D}_{s} +$

MAP Estimation



The MAP



Implications:

- □ The MAP estimator requires to test all the possible supports for the maximization. For the found support, the oracle formula is used.
- □ In typical problems, this process is impossible as there is a combinatorial set of possibilities.
- □ In order to use the MAP approach, we shall turn to approximations.



The MMSE

$$\underline{\hat{\alpha}}^{\mathsf{MMSE}} = \mathsf{E}\left\{\underline{\alpha} \mid \underline{y}\right\} = \sum_{\mathbf{s} \in \Omega} \mathsf{P}(\mathbf{s} \mid \underline{y}) \cdot \mathsf{E}\left\{\underline{\alpha} \mid \underline{y}, \mathbf{s}\right\}$$

This is the oracle for s, as we have seen before $E\{\underline{\alpha} \mid y, s\} = \underline{\hat{\alpha}}_{s} = \mathbf{Q}_{s}^{-1}\underline{h}_{s}$

$$P\left(\underline{s} \middle| \underline{y}\right) \propto \left(\frac{\sigma_{e}^{2}}{\sigma_{\alpha}^{2}}\right)^{\underline{s}} \cdot exp\left\{\frac{\underline{h}_{s}^{\mathsf{T}} \mathbf{Q}_{s}^{-1} \underline{h}_{s}}{2\sigma_{e}^{2}} - \frac{log(det(\mathbf{Q}_{s}))}{2} + \underline{b}^{\mathsf{T}} \underline{s} + \frac{1}{2} \underline{s}^{\mathsf{T}} \mathbf{W} \underline{s}\right\}$$

$$\underline{\hat{\alpha}}^{\mathsf{MMSE}} = \sum_{\mathbf{s}\in\Omega} \mathsf{P}(\mathbf{s} \mid \underline{y}) \cdot \underline{\hat{\alpha}}_{\mathbf{s}}$$



The MMSE

$$\underline{\hat{\alpha}}^{\mathsf{MMSE}} = \mathsf{E}\left\{\underline{\alpha} \mid \underline{y}\right\} = \sum_{\mathbf{s} \in \Omega} \mathsf{P}(\mathbf{s} \mid \underline{y}) \cdot \mathsf{E}\left\{\underline{\alpha} \mid \underline{y}, \mathbf{s}\right\}$$

Implications:

$$\underline{\hat{\alpha}}^{\mathsf{MMSE}} = \sum_{\mathbf{s}\in\Omega} \mathsf{P}(\mathbf{s} \mid \underline{y}) \cdot \underline{\hat{\alpha}}_{\mathbf{s}}$$

□ The best estimator (in terms of L₂ error) is a weighted average of many sparse representations!!!

❑ As in the MAP case, in typical problems one cannot compute this expression, as the summation is over a combinatorial set of possibilities. We should propose approximations here as well.



Part IV Construct Practical Pursuit Algorithms



An OMP-Like Algorithm



The Core Idea – Gather <u>s</u> one element at a time greedily:

- $\hfill\square$ Initialize all entries with {-1}.
- □ Sweep through all entries and check which gives the maximal growth of the expression set it to {+1}.
- □ Proceed this way till the overall expression decreases for any additional {+1}.



A THR-Like Algorithm



The Core Idea – Find the on-support in one-shut:

- \Box Initialize all entries with {-1}.
- □ Sweep through all entries and check for each the obtained growth of the expression sort (down) the elements by this growth.
- □ Choose the first |s| elements as {+1}. Their number should be chosen so as to maximize the MAP expression (by checking 1,2,3, ...).



An Subspace-Pursuit-Like Algorithm



The Core Idea – Add and remove entries from <u>s</u> iteratively:

- □ Operate as in the THR algorithm, and choose the first $|\underline{s}|+\Delta$ entries as $\{+1\}$, where $|\underline{s}|$ is maximizes the posterior.
- □ Consider removing each of the $|\underline{s}| + \Delta$ elements, and see the obtained decrease in the posterior remove the weakest Δ of them.
- \Box Add/remove Δ entries ... similar to [Needell & Tropp `09],[Dai & Milenkovic `09].



An MMSE Approximation Algorithm

$$P\left(\underline{s}|\underline{y}\right) \propto \left(\frac{\sigma_{e}^{2}}{\sigma_{\alpha}^{2}}\right)^{\underline{s}} \cdot \exp\left\{\frac{\underline{h}_{s}^{\mathsf{T}}\mathbf{Q}_{s}^{-1}\underline{h}_{s}}{2\sigma_{e}^{2}} - \frac{\log(\det(\mathbf{Q}_{s}))}{2} + \underline{b}^{\mathsf{T}}\underline{s} + \frac{1}{2}\underline{s}^{\mathsf{T}}\mathbf{W}\underline{s}\right\}$$
$$\underbrace{\hat{\alpha}}^{\mathsf{MMSE}} = \sum_{s \in \Omega} P(s|\underline{y}) \cdot \mathbf{Q}_{s}^{-1}\underline{h}_{s}$$

The Core Idea – Draw random supports from the posterior and average their oracle estimates:

- □ Operate like the OMP-like algorithm, but choose the next +1 by randomly drawing it based on the marginal posteriors.
- □ Repeat this process several times, and average their oracle's solutions.
- □ This is an extension of the Random-OMP algorithm [Elad & Yavneh, `09].



Lets Experiment





Chosen Parameters





Image Denoising & Beyond Via Learned Dictionaries and Sparse representations By: Michael Elad

Gibbs Sampler for Generating s





Results – Representation Error

- □ As expected, the oracle performs the best.
- The Rand-OMP is the second best, being an MMSE approximation.
- MAP-OMP is better than MAP-THR.
- There is a strange behavior in MAP-THR for weak noise – should be studied.
- The plain-OMP is the worst, with practically useless results.





Results – Support Error

- □ Generally the same story.
- Rand-OMP result is not sparse! It is compared by choosing the first |s| dominant elements, where |s| is the number of elements chosen by the MAP-OMP. As can be seen, it is not better than the MAP-OMP.





The Unitary Case

□ What happens when **D** is square and unitary?

□ We know that in such a case, if **W**=0 (i.e., no interactions), then MAP and MMSE have closed-form solutions, in the form of a shrinkage [Turek, Yavneh, Protter, Elad, `10].



□ Can we hope for something similar for the more general **W**? The answer is (partly) positive!



The Posterior Again





Implications

The MAP estimation becomes the following Boolean Quadratic Problem (NP-Hard!)

$$\hat{\underline{s}} = \operatorname{ArgMax}_{s \mid s^{2} = 1} \mathsf{P}\left(\underline{s} \mid \underline{y}\right) = \operatorname{ArgMax}_{s \mid s^{2} = 1} \exp\left\{\underline{q}^{\mathsf{T}} \underline{s} + \frac{1}{2} \underline{s}^{\mathsf{T}} \mathbf{W} \underline{s}\right\}$$



What about exact MMSE? It is possible too! We are working on it NOW.



Results

- ☑ W (of size 64×64) in this experiment has 11 diagonals.
- All algorithms behave as expected.
- Interestingly, the exact-MAP outperforms the approximated MMSE algorithm.
- The plain-OMP (which would have been ideal MAP for the iid case) is performing VERY poorly.





Results

□ As before ...





Part VFirst Steps TowardsAdaptive Recovery



The Grand-Challenge





Maximum-Likelihood?



$$\begin{split} \hat{\mathbf{W}}, \underline{\hat{\mathbf{b}}} &= \operatorname{ArgMax}_{\mathbf{W}, \underline{\mathbf{b}}} \, \mathsf{P}\left(\left\{\underline{\mathbf{s}}_{i}\right\}_{i=1}^{\mathsf{N}} \left| \mathbf{W}, \underline{\mathbf{b}}\right.\right) = \operatorname{ArgMax}_{\mathbf{W}, \underline{\mathbf{b}}} \prod_{i=1}^{\mathsf{N}} \, \mathsf{P}\left(\underline{\mathbf{s}}_{i} \left| \mathbf{W}, \underline{\mathbf{b}}\right.\right) \\ &= \operatorname{ArgMax}_{\mathbf{W}, \underline{\mathbf{b}}} \frac{1}{Z\left(\mathbf{W}, \underline{\mathbf{b}}\right)^{\mathsf{N}}} \prod_{i=1}^{\mathsf{N}} \exp\left\{\underline{\mathbf{b}}^{\mathsf{T}} \underline{\mathbf{s}}_{i} + \frac{1}{2} \underline{\mathbf{s}}_{i}^{\mathsf{T}} \mathbf{W} \underline{\mathbf{s}}_{i}\right\} \end{split}$$

Problem !!! We do not have access to the partition function Z(W,<u>b</u>).



Solution: Maximum-Pseudo-Likelihood

$$\left\{\underbrace{\mathbf{S}}_{i}\right\}_{i=1}^{\mathsf{N}} \longrightarrow \underbrace{\mathsf{MPL-Estimate}}_{\text{of the BM}} \longrightarrow \left\{\mathbf{W}, \underline{b}\right\}$$

$$\overset{\text{Likelihood}}{\text{function}} \mathsf{P}\left(\underline{s} | \mathbf{W}, \underline{b}\right) \longrightarrow \prod_{k=1}^{\mathsf{K}} \mathsf{P}\left(s_{k} | \mathbf{W}, \underline{b}, \underline{s}_{k}^{\mathsf{C}}\right) \xrightarrow{\mathsf{Pseudo-Likelihood}}_{\text{function}}$$

$$\mathsf{P}\left(s_{k} | \underline{s}_{k}^{\mathsf{C}}\right) = \frac{\exp(s_{k}b_{k} + s_{k}\mathbf{W}_{k}^{\mathsf{T}}\underline{s}_{k}^{\mathsf{C}})}{2\cosh(b_{k} + \mathbf{W}_{k}^{\mathsf{T}}\underline{s}_{k}^{\mathsf{C}})} \longrightarrow \left\{\mathbf{W}, \underline{\hat{b}}\right\} = \operatorname*{ArgMax}_{\mathbf{W}, \underline{b}} \left[\begin{array}{c} \operatorname{tr}\left\{\mathbf{S}^{\mathsf{T}}\mathbf{W}\mathbf{S}\right\} + \underline{1}^{\mathsf{T}}\mathbf{S}^{\mathsf{T}}\underline{b}\\ -\underline{1}^{\mathsf{T}}\rho\left(\mathbf{W}\mathbf{S} + \underline{b}\right)\underline{1} \end{array} \right]$$

- □ The function $\rho(\cdot)$ is log(cosh(\cdot)). This is a convex programming task.
- □ MPL is known to be a consistent estimator.
- \Box We solve this problem using SESOP+SD.



Simple Experiment & Results

We perform the following experiment:

- □ W is 64×64 banded, with L=3, and all values fixed (0.2).
- \Box b is a constant vector (-0.1).
- 10,000 supports are created from this BM distribution – the average cardinality is 16.5 (out of 64).
- We employ 50 iterations of the MPL-SESOP to estimate the BM parameters.
- □ Initialization: **W**=0 and <u>b</u> is set by the scalar probabilities for each entry.





Simple Experiment & Results



Testing on Image Patches



□ The results show an improvement of ~0.9dB over plain-OMP.



Part VI Summary and Conclusion



Sparse and Redundant Representation Modeling of Signals – Theory and Applications By: Michael Elad

Today We Have Seen that ...



We intend to explore:

- □ Other pursuit methods
- □ Theoretical study of these algorithms
- \Box K-SVD + this model
- □ Applications ...

What next?

Machine and propose: □ A general set of pursuit methods □ Exact MAP pursuit for

the unitary case

parameters

□ Estimation of the BM

In this work



Sparse and Redundant Representation Modeling of Signals – Theory and Applications By: Michael Elad