

Dictionary Learning Library for MATLAB

Adam S. Charles*

November 9, 2010

1 Introduction

Sparse coding is a concept that is currently obtaining the state of the art results in many fields, including image processing and theoretical neuroscience [3, 5]. The underlying concept of sparse coding is that a high dimensional data can be represented by very few coefficients in some dictionary. Thus although a vector \mathbf{x} may exist in \mathbb{R}^n , it may lie on a manifold of dimension $m \ll n$. \mathbf{x} can thus be written as

$$\mathbf{x} = \mathcal{D}\mathbf{a} + \epsilon \tag{1}$$

where \mathcal{D} is the dictionary (each column is an element of the dictionary), \mathbf{a} is the sparse set of coefficients which generate \mathbf{x} , and ϵ is a white Gaussian noise term indicating imperfections in the generative model.

Two main questions arise in the sparse coding framework with respect to Equation (1):

1. Given a data point \mathbf{x} and a dictionary \mathcal{D} , how can I recover the coefficients \mathbf{a} ?
2. Given a set of data points $\{\mathbf{x}_k\}_{k \in [1, K]}$, can I learn both the coefficients that generated them as well as the dictionary \mathcal{D} in which the decompositions are sparse?

The answer to the first of these questions is well studied in the field of Compressed Sensing, with many readily available packages available [2, 4]. The answer to the second question, however, comes in two distinct flavors, each stemming from a slightly different interpretation of Equation (1). In one interpretation, the Compressed Sensing approach is taken and the answer comes in the form of the K-SVD algorithm [1]. In a different interpretation, Equation (1) is viewed as a probabilistic model, where a MAP estimation is used to determine \mathbf{a} . In this case, a highly kurtotic prior is placed on \mathbf{a} , and the dictionary can be learned in a statistical manner, by sampling from a large number of data samples [6]. This document is for the dictionary learning library version 1.0, which performs dictionary learning over a dataset using the statistical unsupervised method in [6].

*Adam Charles is at the Georgia Institute of Technology, as part of the ECE department and the Neurolaboratory. For any questions/comments, please email at acharles6@gatech.edu

2 Functionality

2.1 Theory of Dictionary Learning

The statistical method of learning dictionaries for sparse coding is based on maximizing the probability distribution over \mathbf{x} , or equivalently minimizing an energy function. In this method, Equation (1) is used to write a likelihood probability distribution on \mathbf{x} ,

$$p(\mathbf{x}|\mathbf{a}, \mathcal{D}) = Z e^{-\frac{\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2}{2\sigma_\epsilon^2}} \quad (2)$$

where σ_ϵ is the variance of the noise term ϵ , and Z is the normalizing constant for the distribution. Now a prior is placed on the coefficients \mathbf{a} . The prior distribution is necessarily highly kurtotic, and while [6] uses a Cauchy distribution the most used currently is the Laplacian distribution. This distribution is chosen because of its relationship to the ℓ^1 norm used in compressed sensing coefficient recovery. The posterior distribution is then

$$p(\mathbf{a}|\mathbf{x}, \mathcal{D}) \propto p(\mathbf{x}|\mathbf{a}, \mathcal{D})p(\mathbf{a}|\mathcal{D}) \quad (3)$$

$$\propto e^{-\frac{\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2}{2\sigma_\epsilon^2}} e^{-\frac{\sqrt{2}}{\sigma_a}\|\mathbf{a}\|_1} \quad (4)$$

where σ_a^2 is the variance of the Laplacian distribution on the coefficients. In Equations (3) and (4), the constant scaling factors are dropped since they do not effect the arg max of the posterior distribution. The conditional in the prior distribution is purely to maintain the concept that this calculation is valid given a dictionary, since technically the prior is independent of \mathcal{D} . Solving the MAP inference problem yields the coefficients \mathbf{a} , given a dictionary \mathcal{D} :

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} (p(\mathbf{a}|\mathbf{x}, \mathcal{D})) \quad (5)$$

$$= \arg \max_{\mathbf{a}} \left(e^{-\frac{\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2}{2\sigma_\epsilon^2}} e^{-\frac{\sqrt{2}}{\sigma_a}\|\mathbf{a}\|_1} \right) \quad (6)$$

$$= \arg \min_{\mathbf{a}} (\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2 + \lambda\|\mathbf{a}\|_1) \quad (7)$$

where $\lambda = 2\sqrt{2}\sigma_\epsilon^2/\sigma_a$. Minimizing the cost function in Equation (7) (the negative log of the posterior) with respect to \mathbf{a} is the optimization given \mathcal{D} , and is a well studied convex optimization problem. In dictionary learning, however, the same energy function needs to be minimized with respect to \mathcal{D} as well.

Finding \mathcal{D} can be viewed as either a maximum likelihood (ML) estimate or another MAP estimate. In the ML version, the optimization is

$$\arg \min_{\mathcal{D}} p(\mathbf{x}|\mathcal{D}) = \arg \min_{\mathcal{D}} \int_{\mathbb{R}^p} p(\mathbf{x}|\mathbf{a}, \mathcal{D})p(\mathbf{a})d\mathbf{a} \quad (8)$$

Which required sampling from the posterior. In [6], Olshausen and Field show that the distribution is tight about the maximum peak $\hat{\mathbf{a}}$, thus the integral can be estimated by finding the maximum and the optimization becomes:

$$\int_{\mathbb{R}^p} p(\mathbf{x}|\mathbf{a}, \mathcal{D}) d\mathbf{a} \approx \langle p(\mathbf{x}|\mathcal{D}, \hat{\mathbf{a}}) \rangle \quad (9)$$

To minimize this likelihood, a gradient descent algorithm can be used with steps on the i^{th} dictionary element (columns of \mathcal{D}) given by

$$\Delta \mathcal{D}_i \propto \langle \hat{\mathbf{a}}(\mathbf{x} - \mathcal{D}\hat{\mathbf{a}}) \rangle \quad (10)$$

where $\langle \rangle$ denotes an average over some sample set of the data.

In order to minimize the likelihood, the following algorithm was set up:

1. Initialize \mathcal{D} to some random dictionary
2. pick a random subsample of the data
3. Find \mathbf{a} given \mathcal{D} for each \mathbf{x} chosen
4. Take a gradient step on \mathcal{D} given by Equation (10)
5. Repeat steps 2 - 4 until convergence

The other method to optimize a dictionary places a prior over the elements of \mathcal{D} as well as \mathbf{a} , resulting in the expended posterior

$$p(\mathbf{a}, \mathcal{D}|\mathbf{x}) = p(\mathbf{x}|\mathbf{a}, \mathcal{D})p(\mathbf{a}|\mathcal{D})p(\mathcal{D}) \quad (11)$$

$$= e^{-\frac{\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2}{2\sigma_\epsilon^2}} e^{-\frac{\sqrt{2}}{\sigma_a}\|\mathbf{a}\|_1} e^{-\frac{\|\mathcal{D}\|_F}{2\sigma_{\mathcal{D}}^2}} \quad (12)$$

In this case the elements of \mathcal{D} are *i.i.d.* Gaussian random variables with zero mean and variance $\sigma_{\mathcal{D}}^2$. The MAP inference is then a joint inference problem given by

$$\{\hat{\mathcal{D}}, \hat{\mathbf{a}}\} = \arg \min_{\{\mathcal{D}, \mathbf{a}\}} (\|\mathbf{x} - \mathcal{D}\mathbf{a}\|_2^2 + \lambda\|\mathbf{a}\|_1 + \lambda_2\|\mathcal{D}\|_F) \quad (13)$$

where $\lambda_2 = \sigma_\epsilon^2/\sigma_{\mathcal{D}}^2$. The learning procedure (in performing an alternating minimization as before) is essentially the same, but with an extra term to the gradient descent step:

$$\Delta \mathcal{D}_i \propto \langle \hat{\mathbf{a}}(\mathbf{x} - \mathcal{D}\hat{\mathbf{a}}) - 2\lambda_2 \mathcal{D}_i \rangle \quad (14)$$

2.2 Code Functionality

The code in the dictionary learning library includes code to learn a dictionary, \mathcal{D} , for data in the form of vectors, image patches, or data cubes (e.g. videos) using either a likelihood or MAP estimate on \mathcal{D} . The code is set to be able to use a preset conjugate gradient descent algorithm to infer the coefficients, but any inference algorithm can be used, with the appropriate wrapper (see Section 2.2.2). Some wrappers are included for some popular inference methods, such as `l1_ls` [4]. For general help with a specific function, type `help` then the function name for comments on its use.

Table 1: Functions included in the Dictionary Learning Library

Function Name	Description
<code>dictionary_learn_script</code>	Example script for setting up dictionary learning on natural images
<code>learn_dictionary</code>	Main function to learn a dictionary, using parfor parallelization
<code>learn_dictionary_spmd</code>	Main function to learn a dictionary, using spmd parallelization
<code>initialize_dictionary</code>	Function to initialize a dictionary
<code>gen_multi_infer</code>	Function to infer coefficients for a data sampling in parallel
<code>dictionary_update</code>	Function to update a dictionary using either a ML or MAP method
<code>cg_l2l1_wrapper</code>	Wrapper for cg l2l1 (included)
<code>l1ls_wrapper</code>	Wrapper for l1 ls (http://www.stanford.edu/~boyd/l1_ls/)
<code>l1ls_nneg_wrapper</code>	Wrapper for l1 ls with non-negativity constraints (in l1 ls package)
<code>SolveMP_wrapper</code>	Wrapper for Matching Pursuit (MP) (http://sparselab.stanford.edu)
<code>SolveOMP_wrapper</code>	Wrapper for Optimized Orthogonal MP (http://sparselab.stanford.edu)
<code>perform_omp_wrapper</code>	Wrapper for regular OMP (http://www.personal.soton.ac.uk/tb1m08/sparsify/sparsify.html)
<code>greed_omp_qr_wrapper</code>	Wrapper for QR-OMP (http://www.personal.soton.ac.uk/tb1m08/sparsify/sparsify.html)
<code>cg_l2l1</code>	Function to infer sparse coefficients using conjugate gradient descent
<code>mintotal</code>	Function used in cg l2l1 for the optimization
<code>basis2image2</code>	Function to reshape a 2D dictionary into a viewable format for plotting
<code>dict_plot1d</code>	Function to plot a vector type dictionary

2.2.1 Included Functions

The included functions are shown in Table 1. The main functions are `learn_dictionary` and `learn_dictionary_spmd`. These functions are equivalent in terms of the input/output characteristics, with the difference being the change in the parallelization scheme.

In `learn_dictionary`, the interior loop where the coefficients of each selected data point is inferred is parallelized using a parfor loop. In `learn_dictionary_spmd`, the parallelization is accomplished by using single program multiple data (spmd) parallelization. The main difference is that parfor apportions the data to the workers as it is selected, while spmd separates the whole data set at the start of the program, and sends only commands to the workers. Different parallelization schemes will be more efficient based on the computer type, data dimensionality and dictionary size. The script `dictionary_learn_script` is an example of starting a MATLAB worker pool, loading data, initializing a dictionary and learning the dictionary for a set of natural images. The file `IMAGES.mat` referenced can be found on Dr. Olshausens website: <http://redwood.berkeley.edu/bruno/sparsenet/>. The full list of options supported by the included code is:

- `opts.data_type` - Type of data ('vector', 'square' or 'cube')
- `opts.grad_type` - Type of cost function for gradient descent: 'norm' or 'forb'.
- `opts.sparse_type` - Type of inference to use. For full set of options, see `multi_infer.m`
- `opts.n_elem` - Number of dictionary elements

- `opts.in_iter` - Number of samples per iteration
- `opts.iters` - Number of iterations to run the algorithm for
- `opts.GD_iters` - Number of gradient descent steps per iteration
- `opts.verb` - 1 if verbose outputs are desired
- `opts.bl_size` - Block size for 'square or 'cube data types
- `opts.dep_size` - Depth size for 'cube data type
- `opts.nneg_dict` - Choose nonnegative values only for the dictionary
- `opts.step_size` - Initial step size for gradient descent
- `opts.decay` - Rate of decay for gradient descent step size
- `opts.lambda` - Lambda value for l1-regulated inference schemes
- `opts.lambda2` - Second lambda value for 'forb energy function
- `opts.tol` - Tolerance for inference schemes
- `opts.ssim_flag` - Choose whether to normalize the data pre-inference
- `opts.std_min` - Minimum standard deviation (use with `.ssim_flag`)
- `opts.save_every` - Number of iterations to save after
- `opts.save_name` - Filename to save the dictionary in
- `opts.bshow` - 0 to not plot intermediary dictionaries. Else plot every `opts.bshow` iterations.
- `opts.disp_size` - Dimensions of figure to display the dictionary

Most of these options have default values and do not need to be set. The bare minimum for running the dictionary learning function is the type of data, the number of dictionary elements, the block size ('square or 'cube data) and depth ('cube data). This dictionary learning package is set up to allow for any inference scheme to be used to find the sparse coefficients at each iteration. All that needs to be used is an appropriate wrapper. A number of wrappers for useful inference schemes have been provided, with the functions available at the corresponding websites in Table 1. For more specific details on the other functions, simply use the help command in MATLAB.

2.2.2 Writing Your Own Wrapper

While wrappers are included for some readily available inference functions, any function can be used as long as a wrapper is written for it and passed to the main function. To write a wrapper, simply use the same inputs/outputs as the existing wrappers and any extra parameters needed can be passed through using the `opts` struct. The inputs to the wrapper must be of the form (`dictionary_n`, `x_im`, `opts`), and the output must be a single output: `coef_vals`. The wrapper simply extracts the necessary options from `opts` and organizes the inputs into the inference function.

References

- [1] M. Aharon, M. Elad, A. Bruckstein, and Y. Katz. K-svd: An algorithm for designing of overcomplete dictionaries for sparse representations. *IEEE Proceedings - Special Issue on Applications of Compressive Sensing & Sparse Representation*.
- [2] S. S. Chen, D. L. Donoho, and M. A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.
- [3] M. Elad, M.A.T. Figueiredo, and Y. Ma. On the role of sparse and redundant representations in image processing. *IEEE Proceedings - Special Issue on Applications of Compressive Sensing & Sparse Representation*, Oct 2008.
- [4] S.-J. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. An interior-point method for large scale l1-regularized least squares. *IEEE Journal on Selected Topics in Signal Processing*, 1(4):606–617, Dec 2007.
- [5] B. A. Olshausen and D.A. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(13):607–609, June 1996.
- [6] B. A. Olshausen and D.A. Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision Research*, 37(23):3311–3325, 1997.