

$$\begin{aligned}
& \underset{\alpha, \mathbf{z}, \mathbf{y}}{\text{minimize}} && \|\mathbf{x} - \mathbf{M}\mathbf{y}\|_2^2 + \sum_{i=1}^L \lambda_i \|\mathbf{z}_i\|_1 \\
& \text{subject to} && \mathbf{D}_1 \alpha_1 = \Phi \mathbf{y} \\
& && \mathbf{P}_{\ell-1} \alpha_{\ell-1} = \mathbf{D}_\ell \alpha_\ell \quad \forall \ell \in \{2, \dots, L\} \\
& && \alpha_\ell = \Phi \mathbf{z}_\ell \quad \forall \ell \in \{1, \dots, L\}
\end{aligned} \tag{1}$$

where Φ is the Fourier transform, $\{\mathbf{P}_i\}_{i \in \{1, \dots, L-1\}}$ are pooling operations, and $\{\mathbf{D}_i\}_{i \in \{1, \dots, L\}}$ are the dictionaries in frequency domain.

$$\begin{aligned}
& \underset{\alpha, \mathbf{D}}{\text{minimize}} && \|\mathbf{x} - \mathbf{M}\mathbf{D}_1 \alpha_1\|_2^2 + \sum_{i=1}^L \lambda_i \|\alpha_i\|_1 \\
& \text{subject to} && \mathbf{P}_{\ell-1} \alpha_{\ell-1} = \mathbf{D}_\ell \alpha_\ell \quad \forall \ell \in \{2, \dots, L\} \\
& && \Phi^{-1} \mathbf{d}_\ell[:, k_\ell] = \mathbf{S}_i \Phi^{-1} \mathbf{d}_\ell[:, k_\ell] \quad \forall \ell \in \{1, \dots, L\}, k_\ell \in \{1, \dots, K_\ell\} \\
& && \frac{1}{P_{\ell-1} K_{\ell-1}} \|\mathbf{d}_\ell[:, k_\ell]\|_2 \leq 1 \quad \forall \ell \in \{1, \dots, L\}, k_{\ell-1} \in \{1, \dots, K_{\ell-1}\}
\end{aligned} \tag{2}$$

where \mathbf{S}_i is a projection to the predened spatial support for \mathbf{d}_i , P_ℓ is the size of each channel at layer ℓ , K_ℓ is the number of channels at layer ℓ (K_0 is the number of input channels), and L is the number of layers.

$$\begin{aligned}
& \underset{\alpha, \mathbf{D}}{\text{minimize}} && \sum_{t=1}^T (\|\mathbf{x}^t - \mathbf{M}\mathbf{D}_1 \alpha_1^t\|_2^2 + \sum_{\ell=1}^L \lambda_\ell \|\alpha_\ell^t\|_1) \\
& \text{subject to} && \mathbf{P}_{\ell-1} \alpha_{\ell-1}^t = \mathbf{D}_\ell \alpha_\ell^t \\
& && \Phi^{-1} \mathbf{d}_\ell[:, k_\ell] = \mathbf{S}_i \Phi^{-1} \mathbf{d}_\ell[:, \ell] \\
& && \frac{1}{P_{\ell-1} K_{\ell-1}} \|\mathbf{d}_\ell[:, \ell]\|_2 \leq 1
\end{aligned} \tag{3}$$

where \mathbf{S}_i is a projection to the predened spatial support for $\mathbf{d}_i[:, k_\ell]$, P_ℓ is the size of each channel at layer ℓ , K_ℓ is the number of channels at layer ℓ (K_0 is the number input channels), and L is the number of layers.

Nothing in this approach precludes us from using filters of different size within the same layer, filter size being controlled by support-projection matrix \mathbf{S}_i . I'm not sure how to word this. Each channel has an effective filter size. If we want to keep this property while introducing varying filter sizes within layers, we will need to adjust the support of the filters of the subsequent layer as well. The subfilters applied to each channel must be adjusted to account for the varying filter sizes from previous layer.

That is, an element of \mathbf{z}_i has a region of support in \mathbf{y} . If different channels have different regions of support, we need to correct for that when we linearly combine filtered channels. Put another way, were we to remove all truncation of convolutions, we need to produce a linear combination of images that are all

of the same size, without zero-padding. I'll derive the math for this later. Feel like this should be a section of the paper, even if we don't implement (which we should, even if we don't fully test on real data). Note, I describe the channel size "without zero padding", but for implementation purposes, zero-padding is the way to go here. It's just that if a filter is s pixels wider, than the subfilters applied to that channel should be Ms pixels thinner in the dictionary for the previous layer.

$$\arg \min_{\alpha} \frac{\rho}{2} \sum_{t=1}^T \|D_1 \alpha_1^t - \Phi y^t + \mu_y^t\|_2^2 + \sum_{\ell=1}^L \|\alpha_\ell^t - \Phi z_\ell^t + \mu_{\alpha_\ell}^t\|_2^2 + \sum_{\ell=2}^L \|P_{\ell-1} \alpha_{\ell-1}^t - D_\ell \alpha_\ell^t + \mu_{D\alpha_\ell}^t\|_2^2 \quad (4)$$

$$\frac{\partial \mathcal{L}(\alpha, \mathbf{z}, \mathbf{y}, \mu)}{\partial \alpha_1^t} = \frac{\rho}{2} (D_1^H D_1 \alpha_1^t + D_1^H (\mu_y^t - \Phi y^t) + \alpha_1^t + (\mu_{\alpha_1}^t - \Phi z_1^t) + P_1^H P_1 \alpha_1^t + \mu_{D\alpha_2}^t - D_2 \alpha_2^t) \quad (5)$$

$$\frac{\partial \mathcal{L}(\alpha, \mathbf{z}, \mathbf{y}, \mu)}{\partial \alpha_1^t} = 0 \implies (\mathbf{I} + P_1^H P_1 + D_1^H D_1) \alpha_1^t - P_1^H D_2 \alpha_2^t = D_1^H (\Phi y^t - \mu_y^t) + \Phi z_1^t - \mu_{\alpha_1}^t - P_1^H \mu_{D\alpha_2}^t \quad (6)$$

$$\frac{\partial \mathcal{L}(\alpha, \mathbf{z}, \mathbf{y}, \mu)}{\partial \alpha_i} = \frac{\rho}{2} (\alpha_i^t - \Phi z_i^t + \mu_{\alpha_i}^t + P_i^H P_i \alpha_i^t - P_i^H (D_{i+1} \alpha_{i+1}^t - \mu_{D\alpha_i}^t) + D_i^H D_i \alpha_i^t - D_i^H (P_{i-1} \alpha_{i-1}^t + \mu_{D\alpha_{i-1}}^t)) \quad (7)$$

$$\frac{\partial \mathcal{L}(\alpha, \mathbf{z}, \mathbf{y}, \mu)}{\partial \alpha_i^t} = 0 \implies -D_i^H P_{i-1} \alpha_{i-1}^t + (\mathbf{I} + P_i^H P_i + D_i^H D_i) \alpha_i^t - P_i^H D_{i+1} \alpha_{i+1}^t = \Phi z_i^t - \mu_{\alpha_i}^t - P_i \mu_{D\alpha_i}^t + D_i^H \mu_{D\alpha_{i-1}}^t \quad (8)$$

Convolutional Sparse Coding Pursuit: Sparse coding pursuit attempts to reconstruct an input signal from a linear combination of atoms. The atoms are selected from a given dictionary, and the coefficients are required to be sparse:

$$\underset{\alpha}{\text{minimize}} \quad \|\mathbf{x} - D\alpha\|_2^2 \quad (9)$$

$$\text{subject to } \|\alpha\|_0 \leq \tau \quad (10)$$

The problem as shown above is NP-hard to solve, but approximate methods have proven sufficient in many contexts. For the purposes of this paper, we will be most interested in the convex relaxation of the above equation (an L-1 constraint is equivalent):

$$\underset{\alpha}{\text{minimize}} \quad \|\mathbf{x} - D\alpha\|_2^2 + \lambda \|\alpha\|_1 \quad (11)$$

In convolutional sparse coding pursuit, the dictionary atoms are convolved with the coefficients to reconstruct the signal.

Ok, this transition isn't going to make sense here.

For signal \mathbf{s} and K dictionary atoms d_1, \dots, d_K , we can formulate a convolutional sparse coding pursuit problem as

$$\underset{\boldsymbol{\alpha}}{\text{minimize}} \left\| \mathbf{x} - \mathbf{M} \sum_{i=1}^K \mathbf{d}_i * \boldsymbol{\alpha}_i \right\|^2 + \lambda \sum_{i=1}^K \|\boldsymbol{\alpha}_i\|_1 \quad (12)$$

where \mathbf{M} is a projection matrix to remove boundary effects and $*$ denotes convolution.

We can rewrite this as

$$\mathbf{D} = [\mathbf{D}_1 \quad \dots \quad \mathbf{D}_K] \quad (13)$$

$$\boldsymbol{\alpha} = \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \vdots \\ \boldsymbol{\alpha}_K \end{bmatrix} \quad (14)$$

$$\underset{\boldsymbol{\alpha}}{\text{minimize}} \left\| \mathbf{x} - \mathbf{M} \mathbf{D} \boldsymbol{\alpha} \right\|_2^2 + \lambda \|\boldsymbol{\alpha}\|_1 \quad (15)$$

Reformulating this as an ADMM problem, we have

$$\underset{\boldsymbol{\alpha}, \mathbf{y}, \mathbf{z}}{\text{minimize}} \quad \left\| \mathbf{x} - \mathbf{M} \mathbf{y} \right\|_2^2 + \lambda \sum_i \|\mathbf{z}_i\|_1 \quad (16)$$

$$\text{subject to} \quad \mathbf{y} = \mathbf{D} \boldsymbol{\alpha} \quad (17)$$

$$\boldsymbol{\alpha}_i = \mathbf{z}_i \quad (18)$$

This will be another abrupt transition, but I need to get this down.

It is necessary to store inverse matrices. Storing inverse matrices for pursuit requires $\mathcal{O}(\sum_{\ell=1}^L K_\ell^2 M N)$.

Storing inverse matrices for updating dictionaries poses a larger challenge. If we learn the dictionary in frequency domain, this will leave us with memory requirements on the order of $\mathcal{O}(K_\ell^2 K_{\ell-1}^2 M N)$. However, it is possible to learn the dictionary in the spatial domain. This changes the memory requirement per layer to $\mathcal{O}(K_\ell^2 K_{\ell-1}^2 F_\ell^2)$ where F_ℓ is the size of the filter. Inconveniently, when using this memory trick, we lose the block-diagonal property we had in frequency domain. I tried unsuccessfully to diagonalize the smaller matrix, but I was unsuccessful. For now, it appears I'm stuck with quartic cost for storage and for multiplication.