

$$\begin{array}{ccc}
 \begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} & = & \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \\
 \mathbf{y} \in \mathbb{R}^m & & \mathbf{B} \in \mathbb{R}^{m \times n} \quad \mathbf{h} \in \mathbb{R}^n
 \end{array} \tag{7.47}$$

In the first expression, we have an example of a sparsely parametrized linear regression model. In the second, we have linear regression with a sparse representation  $\mathbf{h}$  of the data  $\mathbf{x}$ . That is,  $\mathbf{h}$  is a function of  $\mathbf{x}$  that, in some sense, represents the information present in  $\mathbf{x}$ , but does so with a sparse vector.

Representational regularization is accomplished by the same sorts of mechanisms that we have used in parameter regularization.

Norm penalty regularization of representations is performed by adding to the loss function  $J$  a norm penalty on the *representation*. This penalty is denoted  $\Omega(\mathbf{h})$ . As before, we denote the regularized loss function by  $\tilde{J}$ :

$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\mathbf{h}) \tag{7.48}$$

where  $\alpha \in [0, \infty)$  weights the relative contribution of the norm penalty term, with larger values of  $\alpha$  corresponding to more regularization.

Just as an  $L^1$  penalty on the parameters induces parameter sparsity, an  $L^1$  penalty on the elements of the representation induces representational sparsity:  $\Omega(\mathbf{h}) = \|\mathbf{h}\|_1 = \sum_i |h_i|$ . Of course, the  $L^1$  penalty is only one choice of penalty that can result in a sparse representation. Others include the penalty derived from a Student- $t$  prior on the representation (Olshausen and Field, 1996; Bergstra, 2011) and KL divergence penalties (Larochelle and Bengio, 2008) that are especially useful for representations with elements constrained to lie on the unit interval. Lee *et al.* (2008) and Goodfellow *et al.* (2009) both provide examples of strategies based on regularizing the average activation across several examples,  $\frac{1}{m} \sum_i \mathbf{h}^{(i)}$ , to be near some target value, such as a vector with .01 for each entry.

Other approaches obtain representational sparsity with a hard constraint on the activation values. For example, **orthogonal matching pursuit** (Pati *et al.*, 1993) encodes an input  $\mathbf{x}$  with the representation  $\mathbf{h}$  that solves the constrained optimization problem

$$\arg \min_{\mathbf{h}, \|\mathbf{h}\|_0 < k} \|\mathbf{x} - \mathbf{W}\mathbf{h}\|^2, \tag{7.49}$$

where  $\|\mathbf{h}\|_0$  is the number of non-zero entries of  $\mathbf{h}$ . This problem can be solved efficiently when  $\mathbf{W}$  is constrained to be orthogonal. This method is often called