

# Deep Unfolded IRLS-ADMM Network for Classification and Sparse Feature Selection

Xian Yang and Yike Guo

**Abstract**—Two main branches of machine learning methods are model-based methods and deep neural networks. Model-based methods can explicitly include prior knowledge into the model at expense of difficulties in inference, while neural networks are featured in their strong predictive power and straightforward inference approach with the lack of model interpretability. To construct models which are entitled with the advantages of both methods and overcome their problems, the deep unfolding strategy has been developed recently. This paper adopts the idea of deep unfolding to construct a classification and feature selection method. The proposed method is based on the sparse classification; and the iterative inference process of the sparse classification is unfolded into a layer-wise structure analogous to a neural network. Thus, the architecture of our network is fully motivated by the sparse classification method. Different from other neural networks which are essentially black-box methods, our deep unfolded network acts as white-box that features selected in the predictive model can be explicitly returned. Experimental results show the both predictive power and feature selection ability of our methods.

**Index Terms**—Deep unfolded neural network, IRLS-ADMM net, sparse classification, feature selection, white-box method.

## I. INTRODUCTION

Deep unfolding is a recently proposed strategy to construct neural-network-like machine learning models that take the advantage of both model-based methods and neural networks [1]. Model-based methods can easily incorporate prior knowledge, such as sparsity, conditional dependency and latent variable structure, into the model. However, the inference of model-based methods are usually not straightforward and can be both computationally and mathematically intractable that approximation methods such as variational approximation and Monte Carlo sampling are required to iteratively infer variables of interest. Neural networks, on the other hand, organise inference into layers which are typically executed in sequence. Although neural networks become the state of art in many applications, they have the problem of model interpretability. Typical neural methods are black-box methods acting closer to mechanisms than problem-level formulations. Therefore, it comes to the idea of bringing the model-based approach to the task of design the architecture of neural networks, which is called deep unfolding.

The basic idea of deep unfolding is with a given model-based method unfolding the iterations of inference

into a layer-wise structure, which forms into a neural network. To increase the flexibility of the model and the ease of training, model parameters across layers are de-coupled. The resulting formula combines the predictive power of neural network with problem-level formulation of model-based methods. The inference is performed with fixed number of layers and back propagation is used to estimate model parameters. In [1], this deep unfolding strategy is used in the domain of speech enhancement, constructing network based on non-negative matrix factorization [2]. The domain knowledge that signals mix linearly is embodied in the model. Deep unfolding has also been applied to multichannel source separation [3]. A multichannel Gaussian mixture model based neural network has been proved to achieve improved performance of separating signals from simultaneous speakers. In [4], deep unfolded network is applied for supervised topic model. It improves the classification accuracy by replacing the maximisation of lower bound of marginal likelihood in variational Bayesian with the minimisation of cross entropy. Another example is, in [5], a Gaussian conditional random field (GCRF) based neural network is constructed for image denoising. The proposed network can explicitly model the input noise variance and handle a range of noise levels. Along with those models, there are many attempts on constructing neural network based on model-based methods that impose sparse constraints. For example, the work in [6] unfolds the iterations of linear sparse regression method into variants of long short-term memory (LSTM) cells with the results showing that the unfolded network has reduced computational burden in inference. The sparse coding based neural network has also been applied to reconstruct sparse signals from noisy and compressive measurements in [7]. The work in [8] design a neural network based on sparse coding to reconstruct high-quality MR images from under-sampled data, which is effective both in reconstruction accuracy and speed.

In this paper, we focus on constructing a model-based network for both classification and feature selection. The model-based method investigated in this paper is sparse Bayesian classification (SBC) [9], [10]. One of the distinguishing characteristics of SBC is that it does not only build a classification model but also returns a set of features with non-zero weights. In the literature, several classification and feature selection methods are based on SBC [11]-[13]. In our previous work [14], we derive an optimisation-centric classification and feature selection method from SBC. The model is constructed through an iterative convex optimisation procedure instead of a one-step closed form calculation. The cost function is cast using hierarchical Bayesian model, where parameters and hyperparameters are

Manuscript received March 22, 2018; revised May 6, 2018.

The authors are with the Data Science Institute, Imperial College London, SW7 2AZ, UK (e-mail: xian.yang08@imperial.ac.uk, y.guo@imperial.ac.uk).

inferred via the convex optimisation procedure [10], [15]. In this paper, we construct neural network based on this method: the parameters in the model-based method are regarded as hidden nodes in the network, while hyperparameters are learned by training the model through the back-propagation method. We adopt iterative reweighted least squares (IRLS) [16] and alternating direction method of multipliers (ADMM) [8] to get the iterative inference steps. IRLS is used to formulate the problem of finding the step direction for Newton's method as a weighted ordinary least squares problem; AMDD is applied to solve the convex optimisation problem by breaking them into smaller pieces, which has proven to be an efficient variable slitting algorithm with convergence guarantee. After we derived formula for iterative updates, we then design deep architecture consisting of multiple stages, each of which corresponds to an iterative inference step. The weight parameters in the model-based classification model are regarded as nodes in the network, while hyperparameters are learned via stochastic gradient descend (SGD) [17]. The following parts of our paper will first introduce the method in detail. Then the experimental results part will show that the proposed deep unfolded network works well in the aspect of both classification power and feature selection.

## II. DEPP IRLS-ADMM NET FOR CLASSIFICATION AND FEATURE SELECTION

### A. The Sparse Bayesian Classification Based Optimisation Method

Suppose we get a set of input vectors  $\{\mathbf{x}_n\}_{n=1}^N$  along with corresponding targets  $\{y_n\}_{n=1}^N$ . We wish to learn the underlying functional mapping which is defined by a parameterised function  $f(\mathbf{x}; \boldsymbol{\beta}) = \sum_{i=1}^{\aleph} \beta_i \phi_i(\mathbf{x})$ , where the output is the linear weighted sum of  $\aleph$  basis functions and  $\boldsymbol{\beta} = [\beta_1, \beta_2, \dots, \beta_{\aleph}]^T$  contains the parameter. Let  $\Phi$  be the  $N \times \aleph$  design matrix with  $\Phi = [\boldsymbol{\phi}(\mathbf{x}_1), \boldsymbol{\phi}(\mathbf{x}_2), \dots, \boldsymbol{\phi}(\mathbf{x}_N)]^T$ , wherein  $\boldsymbol{\phi}(\mathbf{x}_n) = [\phi_1(\mathbf{x}_n), \phi_2(\mathbf{x}_n), \dots, \phi_{\aleph}(\mathbf{x}_n)]$ . Then we can express the mapping function as  $f(\mathbf{x}; \boldsymbol{\beta}) = \Phi \boldsymbol{\beta}$ . The prior distribution of the weights is assumed to follow a zero-mean isotropic Gaussian:

$$\begin{aligned} \mathcal{P}(\boldsymbol{\beta}|\boldsymbol{\gamma}) &= \prod_{i=1}^{\aleph} \mathcal{N}(\beta_i|0, \gamma_i) \\ &= \prod_{i=1}^{\aleph} (2\pi\gamma_i)^{-\frac{1}{2}} \exp\left\{-\frac{\beta_i^2}{2\gamma_i}\right\}, \end{aligned} \quad (1)$$

where

$$\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_{\aleph}] \in \mathbb{R}^{\aleph}, \boldsymbol{\Gamma} = \text{diag}[\boldsymbol{\gamma}]. \quad (2)$$

The likelihood function  $\mathcal{P}(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x})$  is expressed in the form of the logistic regression model [18]. According to the Bayes' rule, maximisation of posterior is equivalent to finding the maximum over  $\boldsymbol{\beta}$  of  $\log \mathcal{P}(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x})\mathcal{P}(\boldsymbol{\beta}|\boldsymbol{\gamma})$ . The hyperparameter  $\boldsymbol{\gamma}$  is updated by maximising the marginal likelihood, which is equivalent to  $\text{argmin}_{\boldsymbol{\gamma}} - \mathcal{P}(\mathbf{y}|\boldsymbol{\gamma}, \mathbf{x}) = \text{argmin}_{\boldsymbol{\gamma}} - \int \mathcal{P}(\mathbf{y}|\boldsymbol{\beta}, \mathbf{x})\mathcal{P}(\boldsymbol{\beta}|\boldsymbol{\gamma})d\boldsymbol{\beta}$ . In [14], we propose a method to jointly estimate  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$  through a common objective function:

$$\begin{aligned} \text{argmin}_{\boldsymbol{\beta}, \boldsymbol{\gamma}} \quad & \sum_{n=1}^N \log [1 + \exp\{-y_n \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n)\}] \\ & + \boldsymbol{\beta}^T \boldsymbol{\Gamma}^{-1} \boldsymbol{\beta} + \log |\boldsymbol{\Gamma}| + \frac{1}{2} \log |\mathbf{H}(\boldsymbol{\beta}^*)|, \end{aligned} \quad (3)$$

where  $|\mathbf{H}(\boldsymbol{\beta}^*)|$  is the Hessian matrix calculated at mode  $\boldsymbol{\beta}^*$ , which is assumed to be obtained through the minimisation step of  $\boldsymbol{\beta}$  in the iterative optimisation process. This objective function can be expressed in the convex-concave form [19] that the standard iterative optimisation procedure is evoked to get its solution.

### B. Deep IRLS-ADMM Net

The approach in [14] carries out an iterative process to update  $\boldsymbol{\beta}^{k+1}$  and hyper parameters. The target function of obtaining  $\boldsymbol{\beta}$  can be easily formulated in the  $\ell_1$  regularization form as:

$$\begin{aligned} \boldsymbol{\beta}^{k+1} &= \text{argmin}_{\boldsymbol{\beta}} \sum_{n=1}^N \log [1 + \exp\{-y_n \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x}_n)\}] \\ &+ \sum_{i=1}^{\aleph} \|w_i^k \cdot \beta_i\|_{\ell_1}, \end{aligned} \quad (4)$$

where  $\mathbf{w}^k$  is related to hyper-parameters. Rather than considering the iterations as an algorithm, we consider unfolding it as a sequence of layers in a neural network-like architecture. The iteration index is now interpreted as an index to the neural network [1]. Our approach tries to optimise  $\boldsymbol{\beta}^{k+1}$  by treating them as nodes of layers in a neural network parameterised by  $\mathbf{w}^k$ . We optimize the parameter  $\mathbf{w}^k$  via the stochastic gradient descent method by minimising the loss  $\sum_{n=1}^N D(y_n, y_n^*)$ , where  $y_n^*$  is the estimated target and  $D$  is the loss function.

#### 1) The IRLS-ADMM algorithm

To construct the network, we need to obtain the updating steps of  $\boldsymbol{\beta}^{k+1}$  in the form of:

$$\boldsymbol{\beta}^{k+1} = f_{\mathbf{w}^k}(\mathbf{X}, \boldsymbol{\beta}^k). \quad (5)$$

There are a lot of work on using model-based inspiration of novel neural networks, especially on sparse coding problems. However, in sparse coding, the target optimisation function is generally in the form of  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_{\ell_2} + g(\mathbf{x})$ , where  $g(\mathbf{x})$  is the penalty term. Optimisation methods, such as ADMM [8] and proximal gradient method [7], can be used to inspire novel neural networks. For the sparse coding problem, the iterative process for updating  $\mathbf{x}$  can be expressed in the closed form. In our case, however, there is no straightforward way to drive function  $f_{\mathbf{w}^k}$  due to the complex form of the first term in Equation 4.

Our idea is to reformulate the target optimisation function in the form of least squares that typical optimisation methods can be used to find the expression of  $f_{\mathbf{w}^k}$ . We adopt the idea in [16] based on the IRLS formulation of logistic regression. IRLS reformulates the problem of finding the step direction for Newton's method as a weighted ordinary least squares problem. In every iteration, Newton's method first finds a step direction by approximating the objective function by the second order Taylor expression at the current point, and optimising the quadratic approximation in closed-form. The step direction in the  $k$ th iteration is:

$$\boldsymbol{\Omega}^k = \boldsymbol{\beta}^k - \mathbf{H}^{-1}(\boldsymbol{\beta}^k)g(\boldsymbol{\beta}^k), \quad (6)$$

where  $\mathbf{g}$  and  $\mathbf{H}$  are gradient and Hessian matrices. Once the step direction is computed, Newton's method computes the next iteration

$$\boldsymbol{\beta}^{k+1} = (1 - t)\boldsymbol{\beta}^k + t\boldsymbol{\Omega}^k \quad (7)$$

by a linear search over the step size parameter  $t$ .

For the optimisation function in the form of unregularized logistic regression, where the penalty term in Equation 4 is omitted, the gradient and Hessian matrices are in the form of:

$$\mathbf{g}(\boldsymbol{\beta}^k) = - \sum_{n=1}^N y_n \boldsymbol{\phi}(\mathbf{x}_n) [1 - \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\}] \quad (8)$$

and

$$\mathbf{H}(\boldsymbol{\beta}^k) = \boldsymbol{\Phi}^\top \text{diag} [y] \boldsymbol{\Lambda}(\boldsymbol{\beta}^k) \text{diag} [y] \boldsymbol{\Phi}, \quad (9)$$

where

$$\boldsymbol{\Lambda}(\boldsymbol{\beta}^k) = \text{diag} (\{\sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\} [1 - \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\}]\}_{n=1}^N) \quad (10)$$

and

$$\sigma\{f\} = \frac{1}{(1 + e^{-f})}. \quad (11)$$

Let  $\bar{\boldsymbol{\Phi}} = \text{diag} [y] \boldsymbol{\Phi}$ , then we get

$$\mathbf{H}(\boldsymbol{\beta}^k) = \bar{\boldsymbol{\Phi}}^\top \boldsymbol{\Lambda}(\boldsymbol{\beta}^k) \bar{\boldsymbol{\Phi}}. \quad (12)$$

If we define a vector  $\mathbf{v}(\boldsymbol{\beta}^k)$ , whose  $n$ th element is

$$v_n(\boldsymbol{\beta}^k) = y_n \boldsymbol{\phi}(\mathbf{x}_n)^\top \boldsymbol{\beta}^k + \frac{1 - \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\}}{\boldsymbol{\Lambda}(\boldsymbol{\beta}^k)_{n,n}}, \quad (13)$$

then  $\mathbf{g}(\boldsymbol{\beta}^k)$  can be defined in the form of

$$\mathbf{g}(\boldsymbol{\beta}^k) = -\bar{\boldsymbol{\Phi}}^\top \boldsymbol{\Lambda}(\boldsymbol{\beta}^k) (\mathbf{v}(\boldsymbol{\beta}^k) - \bar{\boldsymbol{\Phi}} \boldsymbol{\beta}^k). \quad (14)$$

The step direction in Equation 6 can be expressed as:

$$\boldsymbol{\Omega}^k = (\bar{\boldsymbol{\Phi}}^\top \boldsymbol{\Lambda}(\boldsymbol{\beta}^k) \bar{\boldsymbol{\Phi}})^{-1} \bar{\boldsymbol{\Phi}}^\top \boldsymbol{\Lambda}(\boldsymbol{\beta}^k) \mathbf{v}(\boldsymbol{\beta}^k), \quad (15)$$

which can be regarded as the solution to the following weighted least squares problem:

$$\boldsymbol{\Omega}^k = \underset{\boldsymbol{\Omega}}{\text{argmin}} \|\boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \bar{\boldsymbol{\Phi}} \boldsymbol{\Omega} - \boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \mathbf{v}(\boldsymbol{\beta}^k)\|_{\ell_2}. \quad (16)$$

Therefore, the Newton step direction can be computed by solving the least squares problem defined in Equation 16. The term IRLS refers to the fact that in every iteration the least squares problem has a different diagonal weighting matrix  $\boldsymbol{\Lambda}(\boldsymbol{\beta}^k)$  and  $\mathbf{v}(\boldsymbol{\beta}^k)$ . For the case of  $\ell_1$  regularized logistic regression as in Equation 4, we can get our IRLS formulation as

$$\boldsymbol{\Omega}^k = \underset{\boldsymbol{\Omega}}{\text{argmin}} \frac{1}{2} \|\boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \bar{\boldsymbol{\Phi}} \boldsymbol{\Omega} - \boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \mathbf{v}(\boldsymbol{\beta}^k)\|_{\ell_2} + \lambda \|\mathbf{F}^k \boldsymbol{\Omega}\|_{\ell_1}. \quad (17)$$

Please note that  $\mathbf{F}^k$  in Equation 17 is a diagonal matrix giving weights to elements of  $\boldsymbol{\Omega}$ .

Now, we would like to adopt the ADMM method to get  $\boldsymbol{\Omega}^k$  by solving Equation 17. Suppose

$$\mathbf{A}^k = \boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \bar{\boldsymbol{\Phi}} \quad (18)$$

and

$$\mathbf{b}^k = \boldsymbol{\Lambda}(\boldsymbol{\beta}^k)^{\frac{1}{2}} \mathbf{v}(\boldsymbol{\beta}^k). \quad (19)$$

In ADMM form, the problem in Equation 17 can be written as

$$\begin{aligned} & \text{minimize} \quad \frac{1}{2} \|\mathbf{A}^k \boldsymbol{\Omega} - \mathbf{b}^k\|_{\ell_2} + \lambda \|\mathbf{z}\|_{\ell_1} \\ & \text{subject to} \quad \mathbf{F}^k \boldsymbol{\Omega} - \mathbf{z} = 0, \end{aligned} \quad (20)$$

which yields the ADMM algorithm in the  $i$ th iteration [20]:

$$\begin{aligned} \boldsymbol{\Omega}^{k,i+1} & := (\mathbf{A}^{k\top} \mathbf{A}^k + \rho \mathbf{F}^{k\top} \mathbf{F}^k)^{-1} \\ & \quad (\mathbf{A}^{k\top} \mathbf{b}^k + \rho \mathbf{F}^{k\top} (\mathbf{z}^{k,i} - \mathbf{u}^{k,i})) \\ \mathbf{z}^{k,i+1} & := S_{\lambda/\rho} (\mathbf{F}^k \boldsymbol{\Omega}^{k,i+1} + \mathbf{u}^{k,i}) \\ \mathbf{u}^{k,i+1} & := \mathbf{u}^{k,i} + \mathbf{F}^k \boldsymbol{\Omega}^{k,i+1} - \mathbf{z}^{k,i+1}, \end{aligned} \quad (21)$$

where  $S$  is the soft shareholding operator defined as

$$S_\kappa(\alpha) = \begin{cases} \alpha - \kappa, & \alpha > \kappa \\ 0, & |\alpha| \leq \kappa \\ \alpha + \kappa, & \alpha < -\kappa. \end{cases} \quad (22)$$

The whole process is summarized in Algorithm 1.

---

**Algorithm 1** IRLS-ADMM regularized logistic regression

```

1: Set  $\boldsymbol{\beta}^1 = \mathbf{0}$ 
2: for  $k = 1, \dots, k_{max}$  do
3:   Compute  $\mathbf{A}^k$  and  $\mathbf{b}^k$  using Equation 18 and 19.
4:   Set  $\mathbf{z}^{k,1}$ , and  $\mathbf{u}^{k,1}$  to  $\mathbf{0}$ 
5:   for  $i = 1, \dots, i_{max}$  do
6:     Update  $\boldsymbol{\Omega}^{k,i}$ ,  $\mathbf{z}^{k,i}$  and  $\mathbf{u}^{k,i}$  using Equation 21.
7:   end for
8:   Update  $\boldsymbol{\beta}_k$  using Equation 7, where  $t$  is found using a backtracking line-search that minimizes the objective function.
9: end for
    
```

---

## 2) Deep network design

We follow the idea in [8] to map the iterative process in Algorithm 1 to a data flow graph. In the graph, there are 9 types of nodes described as follows:

**Weighting matrix  $\boldsymbol{\Lambda}^k$  construction layer:** Given  $\boldsymbol{\beta}^k$ , the output of this layer is obtained through:

$$\boldsymbol{\Lambda}_{n,n}^k = \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\} [1 - \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\}], \quad (23)$$

where  $\sigma\{f\} = \frac{1}{(1 + e^{-f})}$ . In the initial stage,  $\boldsymbol{\beta}^0 = \mathbf{0}$ .

**Weighting vector  $\mathbf{v}^k$  construction layer:** Given  $\boldsymbol{\beta}^k$ , the output of this layer is obtained through:

$$\begin{aligned} v_n^k & = y_n \boldsymbol{\phi}(\mathbf{x}_n)^\top \boldsymbol{\beta}^k + \frac{1 - \sigma\{y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)\}}{\boldsymbol{\Lambda}_{n,n}^k} \\ & = y_n \boldsymbol{\phi}(\mathbf{x}_n)^\top \boldsymbol{\beta}^k + 1 + e^{-y_n \boldsymbol{\beta}^{k\top} \boldsymbol{\phi}(\mathbf{x}_n)}. \end{aligned} \quad (24)$$

In the initial stage,  $\boldsymbol{\beta}^0 = \mathbf{0}$ .

**Weighting matrix  $\mathbf{A}^k$  construction layer:** Given  $\boldsymbol{\Lambda}^k$ , the output of this layer is obtained through:

$$\mathbf{A}^k = (\boldsymbol{\Lambda}^k)^{\frac{1}{2}} \bar{\boldsymbol{\Phi}}. \quad (25)$$

**Weighting vector  $\mathbf{b}^k$  construction layer:** Given  $\boldsymbol{\Lambda}^k$  and  $\mathbf{v}^k$ , the output of this layer is obtained through:

$$(\boldsymbol{\Lambda}^k)^{\frac{1}{2}} (\mathbf{v}^k). \quad (26)$$

**Feature  $\boldsymbol{\Omega}^{k,i}$  construction layer:** Given  $\mathbf{z}^{k,i-1}$ ,  $\mathbf{u}^{k,i-1}$ ,  $\mathbf{b}^k$  and  $\mathbf{A}^k$ , the output of this layer is obtained through:

$$\begin{aligned} \Omega^{k,i} &= (\mathbf{A}^{k-1\top} \mathbf{A}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} \mathbf{F}^{k,i})^{-1} \\ &(\mathbf{A}^{k-1\top} \mathbf{b}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1})), \end{aligned} \quad (27)$$

where  $\rho^{k,i}$  and  $\mathbf{F}^{k,i}$  are parameters. When  $i = 1$ ,  $\mathbf{z}^{k,0}$  and  $\mathbf{u}^{k,0}$  are initialized to zeros.

**Linear transform layer  $\mathbf{C}^{k,i}$ :** Given  $\Omega^{k,i}$ , the output of this layer is:

$$\mathbf{C}^{k,i} = \mathbf{F}^{k,i} \Omega^{k,i}, \quad (28)$$

where  $\mathbf{F}^{k,i}$  is the parameter.

**Nonlinear transform layer  $\mathbf{z}^{k,i}$ :** Given  $\mathbf{C}^{k,i}$  and  $\mathbf{u}^{k,i-1}$ , the output of this layer is:

$$\mathbf{z}^{k,i} = S_{\lambda/\rho^{k,i}}(\mathbf{C}^{k,i} + \mathbf{u}^{k,i-1}), \quad (29)$$

where  $\rho^{k,i}$  is the parameter. When  $i = 1$ ,  $\mathbf{u}^{k,0}$  is initialized to zeros.

**Linear transform layer  $\mathbf{u}^{k,i}$ :** Given  $\mathbf{C}^{k,i}$ ,  $\mathbf{u}^{k,i-1}$  and  $\mathbf{z}^{k,i}$ , the output of this layer is:

$$\mathbf{u}^{k,i} = \mathbf{u}^{k,i-1} + \mathbf{C}^{k,i} - \mathbf{z}^{k,i}. \quad (30)$$

When  $i = 1$ ,  $\mathbf{u}^{k,0}$  is initialized to zeros.

**Linear transform layer  $\beta^k$ :** Given  $\Omega^{k-1,i_{\max}}$  and  $\beta^{k-1}$ , the output of this layer is defined as

$$\beta^k = (1 - t^k) \beta^{k-1} + t^k \Omega^{k-1,i_{\max}}, \quad (31)$$

where  $t^k$  is the parameter.

An example deep network with  $k_{\max} = 3$  and  $i_{\max} = 4$  are shown in Fig. 1.

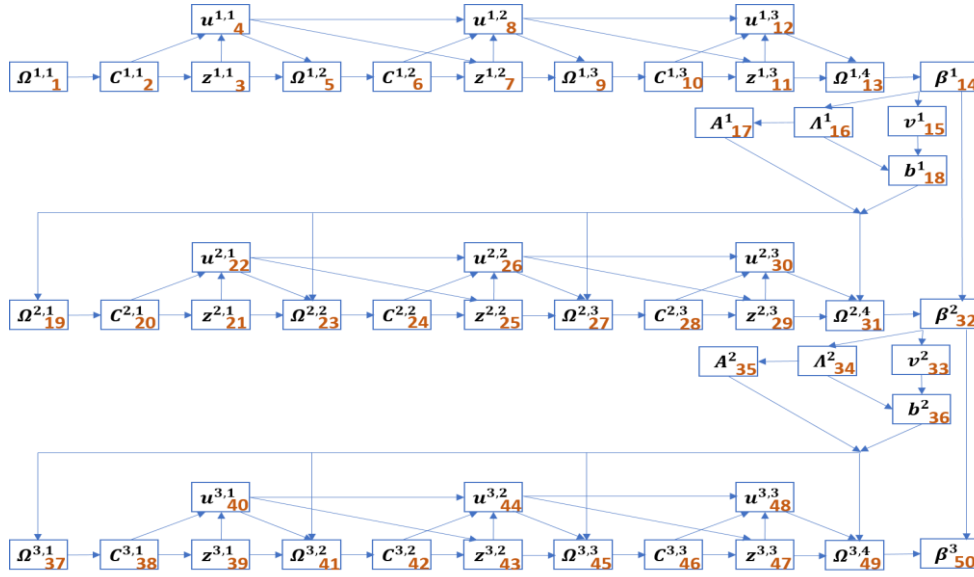


Fig. 1. The architecture of deep IRLS-ADMM net with  $k_{\max} = 3$  and  $i_{\max} = 4$ .

### III. NETWORK TRAINING

#### A. Network Loss

The loss function we choose for network training is:

$$E = \sum_{n=1}^N \log [1 + \exp \{-y_n \beta^\top \phi(x_n)\}], \quad (32)$$

where  $\beta$  is the network output. We learn the parameters  $\Theta = \{\rho^{k,i}, \mathbf{F}^{k,i}, t^k\}$  by minimising the loss w.r.t. them using the stochastic gradient descent method.

#### B. Gradient Computation by Back-Propagation over Data Flow Graph

Fig. 2 shows 19 types of nodes, most of which have multiple inputs and outputs. The parameters we need to learn are  $\Theta = \{\rho^{k,i}, \mathbf{F}^{k,i}, t^k\}$  for  $k \in \{1,2,3\}$  and  $i \in \{1,2,3,4\}$ . Please refer to the Appendix A for gradients computation for each node.

### IV. EXPERIMENTS

#### A. Simulation Results

We generate a data matrix with 500 samples and 50 features from the normal distribution. To check the ability of

finding correlated features, we split the features with non-zero weights into two sets,  $S^1$  and  $S^2$ , each of which contains 4 true features. Then, we initialise a design matrix  $\Phi = \mathbf{x}$  with 500 samples and 50 features from the normal distribution. Let  $\phi_i(\mathbf{x})$  be the  $i$ th column of  $\Phi$ . We set  $\phi_{S_m^2}(\mathbf{x}) = \phi_{S_m^1}(\mathbf{x}) + \mathcal{N}(0,0.1)$ , where  $S_m^1$  and  $S_m^2$  denote the  $m$ th element in each feature set. In this way, we get data generated from correlated features. The target variable vector  $\{y_n\}_{n=1}^{500}$  is generated by the linear model using  $\beta$  and  $\Phi$  with additive independent identical distributed Gaussian noises, where the standard deviation of noise varies ranging from  $\{0,0.1,0.5,1\}$ . In our method,  $\rho^{k,i}$  and  $t^k$  are initialized to be 1. Each element in the diagonal matrix  $\mathbf{F}^{k,i}$  is initialized to be 10.

We will first use the whole dataset with no additive noise to construct the predictive model, from which we compare the estimated weights  $\beta$  from our method with the regularized logistic regression method (ADMM). It is to investigate the ability of selecting real features using our method. Moreover, we would like to compare the scaled magnitudes of estimated weights  $\beta$  with real weights and estimated weights from the regularized logistic regression method. We scale the estimated results to make the first real nonzero feature have the same value. In this way, we can intuitively see whether the estimation has maintained the

ratio in magnitudes of different weights. Thus the rank of the importance of features, which is usually determined by the magnitudes, can be obtained correctly. This is quite important when we use our method as a feature selection approach to select top relevant features in a cross validation process. The results are shown in Fig. 3. Fig. 3 shows the weights from our method and ADMM, from which we can see that there are 50 features, out of which 8 features are

nonzero. Our method has successfully detected these nonzero ones with magnitudes significantly larger than the others, while the ADMM method fails to detect 3 nonzero feature. These observations are quite encouraging, showing that our method can be used a good tool to select relevant features. Most classification methods are not guaranteed to have this characteristic.

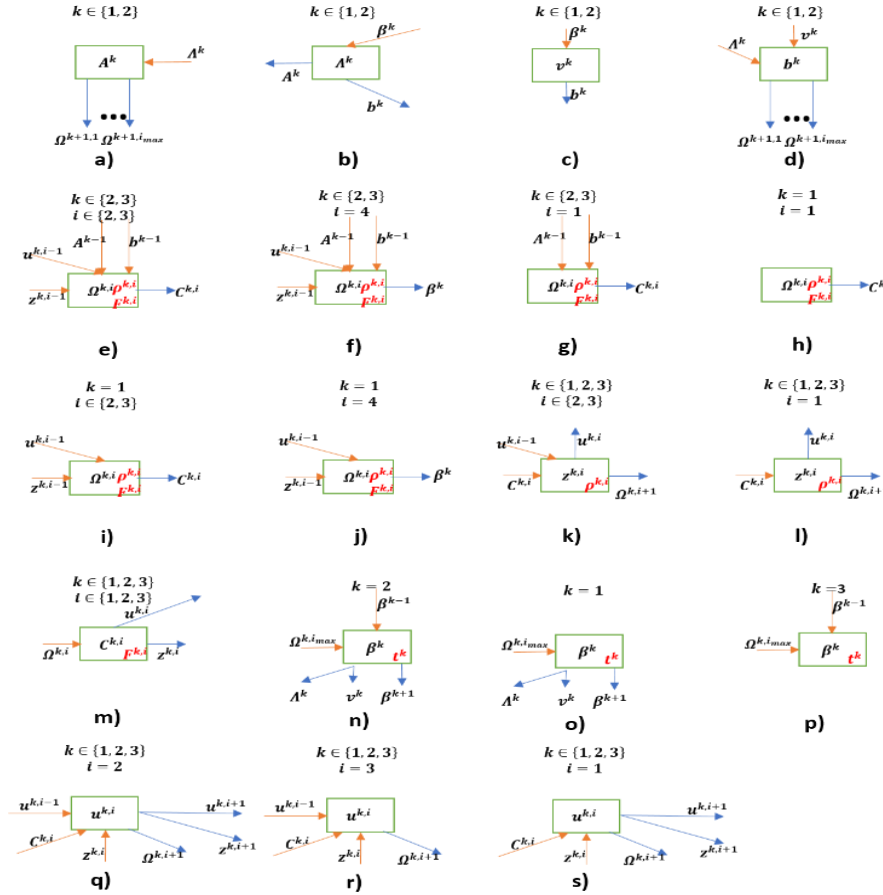


Fig. 2. Illustration of 19 types of nodes. All solid arrows indicate data flow in forward pass.

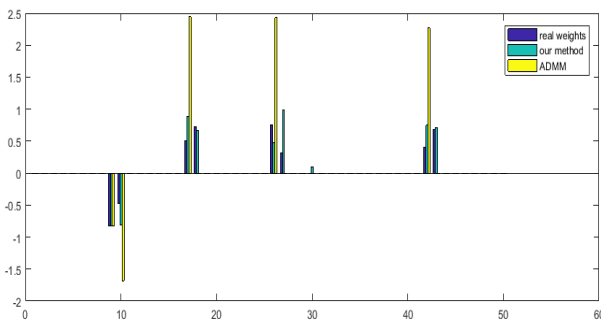


Fig. 3. Comparison of weights estimated from our method and ADMM.

We then investigate the performance of our method with different levels of noise. The performance, including prediction accuracy, false positive (FP) and false negative (FN) rates, is evaluated by the 5-fold cross validation (CV). To calculate FP and FN, we need to compare the selected features across CV with real non-zero features. As in the cross validation process, different training dataset is used for feature selection and predictive model construction in each fold, selected feature sets from all folds may vary due to the variation of training datasets. To select a feature set which is

stable with small fluctuations of the input dataset and also has good predictive accuracy, we use the method from [21]. In the  $k$ th fold of cross validation, the whole dataset  $\mathbf{D}$  is split into two subsets: CV training dataset  $\mathbf{D}_k$  and CV testing dataset  $\mathbf{D}_{\setminus k}$ . Our method can work as a feature selection method on the training dataset  $\mathbf{D}_k$  to rank and select the top  $q$  features, labelled as  $\mathbf{V}_{q,k}$ . After features have been selected, our method then constructs a predictive model for classification using  $\mathbf{V}_{q,k}$ . The prediction results at this CV fold are recorded for later evaluation. To get the complete prediction results, we repeat the above steps for all folds of CV. The method presented in [21] returns an optimal feature set  $\tilde{\mathbf{V}}_q$  with an associated performance score  $\tilde{P}_q$  under each value of  $q$ . The score  $\tilde{P}_q$  is calculated according to the 6<sup>th</sup> strategy proposed in [22] to assess the prediction accuracy and stability of features (the details of calculating this score can be found in [22]). By checking the maximum value of  $\tilde{P}_q$ , we can determine the optimal value of  $q$  and the corresponding optimal feature set  $\tilde{\mathbf{V}}_q$ . The detected optimal set can then be used to construct a predictive model for future prediction.

Fig. 4 shows the change of scores  $\tilde{P}_q$  under different settings of feature size  $q$  and noise level. We can see that the scores of our method under different noise levels peak when the value of  $q$  is close to the number of real nonzero features. In contrast to our method, the optimal value of  $q$  for the ADMM method is around 4, which is different from the real number of nonzero features. From Fig. 4, we can expect that our method works better than the ADMM method in the aspect of feature selection. In the following experiments, we would like to check the performance of these methods averaged over 20 randomly sampled datasets. The results are shown in Table I. The feature set size  $q$  is chosen to be the optimal value detected from Fig. 4. To show the ability of detecting real positives (i.e. real nonzero features), we also present the results of false positive (FP) and false negative (FN) rates in Table I. We can see that under different noise levels, the accuracy achieved by different methods are all maintained at high levels. However, the false negative rates from ADMM are much higher than the rates from our method. This is because, the ADMM method can detect correlated features that some real features are ignored.

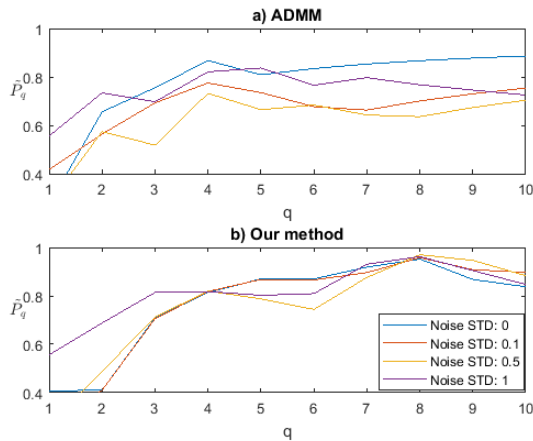


Fig. 4. The scores achieved by different methods with the size of feature set varying from 1 to 8 and standard deviation of noise chosen from 0, 0.1, 0.5 and 1.

TABLE I: THE RESULTS OF OUR METHOD AND ADMM UNDER DIFFERENT NOISE LEVELS. THE CLASSIFICATION ACCURACY, FALSE POSITIVE AND FALSE NEGATIVE RATES ARE COMPARED

Noise STD	Method	Accuracy	FP	FN
0	Our method	0.97	0.02	0
	ADMM	0.98	0.09	0.45
0.1	Our method	0.97	0.02	0
	ADMM	0.98	0.09	0.46
0.5	Our method	0.97	0.02	0
	ADMM	0.96	0.08	0.42
1	Our method	0.96	0.02	0
	ADMM	0.98	0.08	0.44

## B. Real Examples

### 1) Cervical cancer datasets

To demonstrate the applicability of our method in real world problems, we first use the public dataset to test our performance, which can be downloaded from <https://archive.ics.uci.edu/ml/datasets/Cervical+cancer+%28Risk+Factors%29>. This dataset focuses on the prediction of

diagnosis of cervical cancer. There are 4 indicators (Hinselmann, Schiller, Citology and Biopsy) of disease status. We generate a consensus indicator of disease state that if more than 2 indicators saying the patient has the cancer then we say the patient has the cancer. The dataset comprises of demographic information, habits and historic medical records of patients. We choose records of 737 patients with no missing values and 21 features with at least 10 nonzero records across all patients. The detailed descriptions of features in this risk factor dataset can be found in [23].

We first use the whole dataset to select features used in predictive model construction. The estimated weights are shown in Fig. 5. From Fig. 5, we can see that 5 out of 21 features are mainly used for predicting diagnosis of cervical cancer, which are 'Number of sexual partners', 'First sexual intercourse (age)', 'Number of pregnancies', 'Smokes' and 'Hormonal Contraceptives'. Among these features, the 'First sexual intercourse (age)' feature has the largest absolute magnitude of weight. Its negative value shows that people have their first sexual intercourse at early age are more likely to have cervical cancer. These findings can help researchers to get a proper list of risk factors for cervical cancer prediction. After investigating the ability of feature selection, we then also carry out a 5-fold CV process to evaluate the prediction accuracy of our method. The mean accuracy of our method is as high as 0.95.

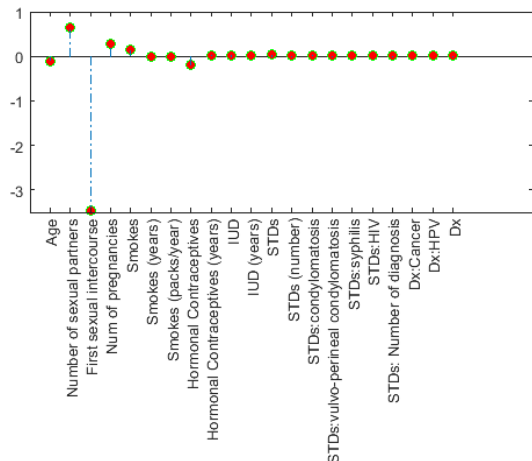


Fig. 5. Weights estimated from our method in the predictive model for cervical cancer risk factors dataset.

### 2) Embryonal tumour gene expression data analysis

We also use a public available gene expression dataset of the central nervous system embryonal tumours from the study in [24]. The raw data can be downloaded from <http://archive.broadinstitute.org/mpr/CNS>. We selected 10 CNS medulloblastomas (MD) samples and 10 non-neuronal origin malignant gliomas (Mglio) samples to show the performance of our method in classifying two tumour types. The samples were hybridised on Affymetrix HuGeneFL GeneChip arrays. We first preprocessed the raw data using GCRMA with empirical Bayes estimate [25]. Then we filtered out probe sets which are either not annotated or have little variability across samples. We select probes whose  $p$  value from  $t$  test for sample type comparison is smaller than 0.05 and fold change is larger than 2. Probes for 350 genes were remained after preprocessing.

Our method can find differences between two tumour types at molecular level. We construct a classifier using the

selected 20 samples with the accuracy of tumour type prediction approaching to 100%. The beauty of our method is that it does not only have strong predictive power, but also selects relevant features that could be candidates of disease biomarkers. Fig. 6a) shows the weights of features in our classification model, where 21 features are mainly used for classification. Features with nonzero weights can be regarded as molecular features distinguishing tumours. By looking at Fig. 6b), we can see that many of these features are highly correlated, telling that our method does not discard features from correlated ones.

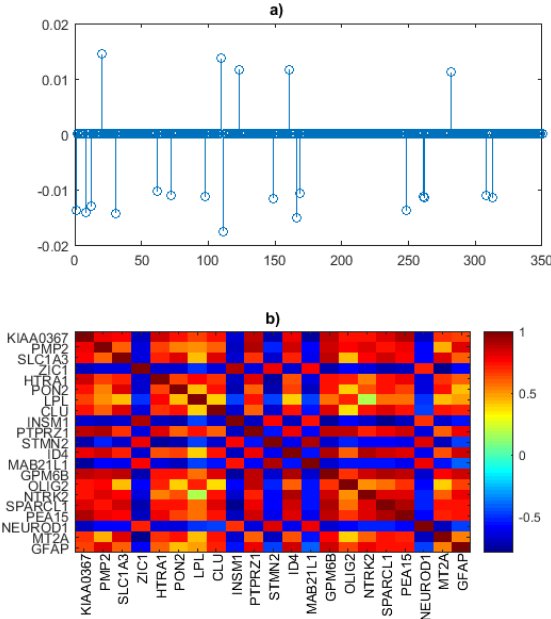


Fig. 6. Results of gene expression analysis: a) weights of features from the classification models using our method; b) Heatmap of the correlation matrix for 21 genes selected by our classifier.

## V. CONCLUSION

This paper adopts the idea of deep unfolding to design a neural network based on IRLS-ADMM for both classification and feature selection. The proposed IRLS-ADMM net is a novel deep neural network. Different from conventional neural networks, which are essentially black-box methods, our method works as a white-box that the network structure are designed according to the iterative updating steps in the model-based model. One advantage of our method is that the network can return a set of features that are used for prediction while maintain good prediction accuracy. Unlike other neural networks, the proposed method can explicitly select features which can be used as a feature selection tool in many applications, such as clinical variable selection and biomarker discovery. Our experimental results have shown that our method has taken the advantages of model-based methods for explicitly incorporating sparsity into model construction and also neural networks with their strong predictive power and straightforward parameter inference process.

### APPENDIX A: GRADIENT COMPUTATION FOR DEEP IRLS-ADMM NET

**Node a):** As shown in Fig. 2a), this node has one input:  $\Lambda^k$

for  $k \in \{1,2\}$ . Its output  $A^k$  is the input to compute  $\Omega^{k,i}$  for all  $i \in \{1,2,3,4\}$ . This layer has no parameters. The loss w.r.t. output is

$$\frac{\partial E}{\partial A^k} = \sum_{i=1}^4 \frac{\partial E}{\partial \Omega^{k+1,i}} \frac{\partial \Omega^{k+1,i}}{\partial A^k}. \quad (33)$$

The operation in this layer is:

$$A^k = (\Lambda^k)^{\frac{1}{2}} \bar{\Phi}. \quad (34)$$

We compute the gradients of the output in this layer w.r.t its input

$$\frac{\partial A^k}{\partial \Lambda_{n,n}^k} = \frac{1}{2} (\Lambda^k)^{-\frac{1}{2}} \bar{I}_{n,n}^k \bar{\Phi}, \quad (35)$$

where  $\bar{I}_{n,n}^k$  is the matrix whose entries are zero except the  $n$ th element on the diagonal equalling to 1.

**Node b):** As shown in Fig. 2b), this node has one input:  $\beta^k$ . Its output  $\Lambda^k$  is the input to compute  $A^k$  and  $b^k$  for  $k \in \{1,2\}$ . This layer has no parameters. The loss w.r.t. output is

$$\frac{\partial E}{\partial \Lambda^k} = \frac{\partial E}{\partial A^k} \frac{\partial A^k}{\partial \Lambda^k} + \frac{\partial E}{\partial b^k} \frac{\partial b^k}{\partial \Lambda^k}. \quad (36)$$

The operation in this layer is:

$$\Lambda_{n,n}^k = \sigma\{y_n \beta^k \top \phi(x_n)\} [1 - \sigma\{y_n \beta^k \top \phi(x_n)\}], \quad (37)$$

where  $\sigma\{f\} = \frac{1}{(1+e^{-f})}$ . We compute the gradients of the output in this layer w.r.t its input, that is  $\frac{\partial \Lambda^k}{\partial \beta_m^k}$  for all  $m$ .  $\frac{\partial \Lambda^k}{\partial \beta_m^k}$  is a diagonal matrix, whose  $n$ th entry on the diagonal is:

$$\begin{aligned} \frac{\partial \Lambda_{n,n}^k}{\partial \beta_m^k} &= y_n \left( \frac{e^{y_n \beta^k \top \phi(x_n)}}{(1 + e^{y_n \beta^k \top \phi(x_n)})^3} \right. \\ &\quad \left. - \frac{e^{-y_n \beta^k \top \phi(x_n)}}{(1 + e^{-y_n \beta^k \top \phi(x_n)})^3} \right) \phi_m(x_n), \end{aligned} \quad (38)$$

where  $\phi_m(x_n)$  is the  $m$ th element of  $\phi(x_n)$ .

**Node c):** As shown in Fig. 2c), this node has one input:  $\beta^k$ . Its output  $v^k$  is the input to compute  $b^k$  for  $k \in \{1,2\}$ . This layer has no parameters. The loss w.r.t. output is

$$\frac{\partial E}{\partial v^k} = \frac{\partial E}{\partial b^k} \frac{\partial b^k}{\partial v^k}. \quad (39)$$

The operation in this layer is:

$$v_n^k = y_n \phi(x_n) \top \beta^k + 1 + e^{-y_n \beta^k \top \phi(x_n)}. \quad (40)$$

We compute the gradients of the output in this layer w.r.t its input, that is  $\frac{\partial v^k}{\partial \beta^k}$ , whose entry in the  $n$ th row and  $m$ th column is

$$\begin{aligned} \frac{\partial v_n^k}{\partial \beta_m^k} &= \frac{\partial y_n \phi(x_n) \top \beta^k + 1 + e^{-y_n \beta^k \top \phi(x_n)}}{\partial \beta_m^k} \\ &= y_n \phi_m(x_n) (1 - e^{-y_n \beta^k \top \phi(x_n)}), \end{aligned} \quad (41)$$

where  $\phi_m(x_n)$  is the  $m$ th element of  $\phi(x_n)$ .

**Node d):** As shown in Fig. 2d), this node has two sets of inputs:  $\Lambda^k$  and  $v^k$ . Its output  $b^k$  is the input to compute  $\Omega^{k,i}$  for all  $i \in \{1,2,3,4\}$  and  $k \in \{1,2\}$ . This layer has no parameters. The loss w.r.t. output is

$$\frac{\partial E}{\partial \mathbf{b}^k} = \sum_{i=1}^4 \frac{\partial E}{\partial \Omega^{k+1,i}} \frac{\partial \Omega^{k+1,i}}{\partial \mathbf{b}^k}. \quad (42)$$

The operation in this layer is:

$$\mathbf{b}^k = (\Lambda^k)^{\frac{1}{2}}(\mathbf{v}^k). \quad (43)$$

We compute the gradients of the output in this layer w.r.t its input

$$\frac{\partial \mathbf{b}^k}{\partial \Lambda_{n,n}^k} = \frac{1}{2} (\Lambda^k)^{-\frac{1}{2}} \bar{\mathbf{I}}_{n,n}^k \mathbf{v}^k \quad (44)$$

and

$$\frac{\partial \mathbf{b}^k}{\partial \mathbf{v}^k} = (\Lambda^k)^{\frac{1}{2}}, \quad (45)$$

where  $\bar{\mathbf{I}}_{n,n}^k$  is the matrix whose entries are zero except the  $n$ th element on the diagonal equalling to 1.

**Node e):** As shown in Fig. 2e), this node has four sets of inputs:  $\mathbf{b}^{k-1}$ ,  $\mathbf{A}^{k-1}$ ,  $\mathbf{z}^{k,i-1}$  and  $\mathbf{u}^{k,i-1}$  for  $k \in \{2,3\}$  and  $i \in \{2,3\}$ . Its output  $\Omega^{k,i}$  is the input to compute  $\mathbf{C}^{k,i}$ . Parameters of this layer are  $\rho^{k,i}$  and  $\mathbf{F}^{k,i}$ . Please note that these two parameters are shared among different layers. Here we just calculate its gradient via the node in this layer. To get the complete the gradient, we need to sum all gradients shared among layers. We will calculate this later. The gradients of loss w.r.t. the parameters can be computed as

$$\frac{\partial E}{\partial \rho^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} \quad (46)$$

and

$$\frac{\partial E}{\partial \mathbf{F}^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \mathbf{F}^{k,i}}, \quad (47)$$

where

$$\frac{\partial E}{\partial \Omega^{k,i}} = \frac{\partial E}{\partial \mathbf{C}^{k,i}} \frac{\partial \mathbf{C}^{k,i}}{\partial \Omega^{k,i}}. \quad (48)$$

To get a complete representation of gradients, we need to do the following calculations. The operation in this layer is:

$$\Omega^{k,i} = (\mathbf{A}^{k-1\top} \mathbf{A}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} \mathbf{F}^{k,i})^{-1} (\mathbf{A}^{k-1\top} \mathbf{b}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1})). \quad (49)$$

Assume

$$\mathbf{Q} = (\mathbf{A}^{k-1\top} \mathbf{A}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} \mathbf{F}^{k,i})^{-1}. \quad (50)$$

The gradients of output in this layer w.r.t. parameters are:

$$\begin{aligned} \frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} &= -Q^2 \mathbf{F}^{k,i\top} \mathbf{F}^{k,i} (\mathbf{A}^{k-1\top} \mathbf{b}^{k-1} \\ &+ \rho^{k,i} \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1})) \\ &+ Q \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1}) \end{aligned} \quad (51)$$

and

$$\begin{aligned} \frac{\partial \Omega^{k,i}}{\partial \mathbf{F}_{m,m}^{k,i}} &= -2\rho^{k,i} Q^2 \bar{\mathbf{F}}_{m,m}^{k,i} (\mathbf{A}^{k-1\top} \mathbf{b}^{k-1} \\ &+ \rho^{k,i} \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1})) \\ &+ \rho^{k,i} Q \bar{\mathbf{I}}_{m,m}^{k,i} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1}), \end{aligned} \quad (52)$$

where  $\bar{\mathbf{F}}_{m,m}^{k,i}$  is the  $m$ th entry on diagonal of  $\mathbf{F}^{k,i}$ ,  $\bar{\mathbf{F}}_{m,m}^{k,i}$  and  $\bar{\mathbf{I}}_{m,m}^{k,i}$  are matrices whose entries are zero except the  $m$ th element on the diagonal equalling to  $\bar{\mathbf{F}}_{m,m}^{k,i}$  and 1, respectively. We also compute the output in this layer w.r.t. its inputs as follows.

$$\frac{\partial \Omega^{k,i}}{\partial \mathbf{b}^{k-1}} = \mathbf{Q} \mathbf{A}^{k-1\top} \quad (53)$$

$$\frac{\partial \Omega^{k,i}}{\partial \mathbf{z}^{k,i-1}} = \rho^{k,i} \mathbf{Q} \mathbf{F}^{k,i\top} \quad (54)$$

$$\frac{\partial \Omega^{k,i}}{\partial \mathbf{u}^{k,i-1}} = -\rho^{k,i} \mathbf{Q} \mathbf{F}^{k,i\top} \quad (55)$$

$$\begin{aligned} \frac{\partial \Omega^{k,i}}{\partial \mathbf{A}_{n_x, n_y}^{k-1}} &= -Q^2 \bar{\mathbf{A}}_{n_x, n_y}^{k-1} (\mathbf{A}^{k-1\top} \mathbf{b}^{k-1} \\ &+ \rho^{k,i} \mathbf{F}^{k,i\top} (\mathbf{z}^{k,i-1} - \mathbf{u}^{k,i-1})) \\ &+ Q \bar{\mathbf{T}}_{n_x, n_y}^{k-1} \mathbf{b}^{k-1}. \end{aligned} \quad (56)$$

We have

$$\bar{\mathbf{A}}_{n_x, n_y}^{k-1} = \tilde{\mathbf{A}}_{n_x, n_y}^{k-1} + \tilde{\mathbf{A}}^{k-1\top}_{n_x, n_y}, \quad (57)$$

where  $\tilde{\mathbf{A}}_{n_x, n_y}^{k-1}$  is a matrix containing zeros except the  $n_y$ th row equals to the  $n_x$ th row of  $\mathbf{A}^{k-1}$ . And  $\tilde{\mathbf{T}}_{n_x, n_y}^{k-1}$  is a matrix containing zeros except the entry in the  $n_x$ th row and  $n_y$ th column equals to 1.

**Node f):** As shown in Fig. 2f), this node is nearly same as the node in Fig. 2e). The only difference is that its output is the input to compute  $\beta_k$ .

$$\frac{\partial E}{\partial \Omega^{k,i}} = \frac{\partial E}{\partial \beta_k} \frac{\partial \beta_k}{\partial \Omega^{k,i}} \quad (58)$$

for  $k \in \{2,3\}$  and  $i = 4$ .

**Node g):** As shown in Fig. 2g), this node is similar with the node in Fig. 2e). The only difference is that it only have two inputs:  $\mathbf{b}^{k-1}$  and  $\mathbf{A}^{k-1}$  for  $k \in \{2,3\}$  and  $i = 1$ . The operation in this layer is:

$$\Omega^{k,i} = (\mathbf{A}^{k-1\top} \mathbf{A}^{k-1} + \rho^{k,i} \mathbf{F}^{k,i\top} \mathbf{F}^{k,i})^{-1} \mathbf{A}^{k-1\top} \mathbf{b}^{k-1}. \quad (59)$$

The gradients of output in this layer w.r.t. parameters are:

$$\frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} = -Q^2 \mathbf{F}^{k,i\top} \mathbf{F}^{k,i} \mathbf{A}^{k-1\top} \mathbf{b}^{k-1} \quad (60)$$

and

$$\frac{\partial \Omega^{k,i}}{\partial \mathbf{F}_{m,m}^{k,i}} = -2\rho^{k,i} Q^2 \bar{\mathbf{F}}_{m,m}^{k,i} \mathbf{A}^{k-1\top} \mathbf{b}^{k-1}. \quad (61)$$

The output in this layer w.r.t. its inputs as follows.

$$\frac{\partial \Omega^{k,i}}{\partial \mathbf{b}^{k-1}} = \mathbf{Q} \mathbf{A}^{k-1\top} \quad (62)$$

and

$$\begin{aligned} \frac{\partial \Omega^{k,i}}{\partial \mathbf{A}_{n_x, n_y}^{k-1}} &= -Q^2 \bar{\mathbf{A}}_{n_x, n_y}^{k-1} \mathbf{A}^{k-1\top} \mathbf{b}^{k-1} \\ &+ Q \bar{\mathbf{T}}_{n_x, n_y}^{k-1} \mathbf{b}^{k-1}. \end{aligned} \quad (63)$$

**Node h):** As shown in Fig. 2h), this node has no input. Its output  $\Omega^{k,i}$  is the input to compute  $\mathbf{C}^{k,i}$  for  $k = 1$  and  $i = 1$ . The operation in this layer is:



$$\Omega^{1,1} = (A^{0\top} A^0 + \rho^{1,1} F^{1,1\top} F^{1,1})^{-1} A^{0\top} b^0. \quad (64)$$

where

$$A^0 = (\Lambda^0)^{\frac{1}{2}} \bar{\Phi} \quad (65)$$

$$b^0 = (\Lambda^0)^{\frac{1}{2}} v^0. \quad (66)$$

As  $\beta^0 = \mathbf{0}$ , we get  $\Lambda^0 = \frac{1}{4}$  and  $v^0 = 2$ . Then we have

$$A^0 = \frac{1}{2} \bar{\Phi} \quad (67)$$

$$b^0 = \vec{1}. \quad (68)$$

Parameters of this layer are  $\rho^{1,1}$  and  $F^{1,1}$ . The gradients of loss w.r.t. the parameters can be computed as

$$\frac{\partial E}{\partial \rho^{1,1}} = \frac{\partial E}{\partial \Omega^{1,1}} \frac{\partial \Omega^{1,1}}{\partial \rho^{1,1}} \quad (69)$$

and

$$\frac{\partial E}{\partial F^{1,1}} = \frac{\partial E}{\partial \Omega^{1,1}} \frac{\partial \Omega^{1,1}}{\partial F^{1,1}}, \quad (70)$$

where

$$\frac{\partial E}{\partial \Omega^{1,1}} = \frac{\partial E}{\partial C^{1,1}} \frac{\partial C^{1,1}}{\partial \Omega^{1,1}}. \quad (71)$$

The gradients of output in this layer w.r.t. parameters are:

$$\frac{\partial \Omega^{1,1}}{\partial \rho^{1,1}} = -(A^{0\top} A^0 + \rho^{1,1} F^{1,1\top} F^{1,1})^{-2} F^{1,1\top} F^{1,1} A^{0\top} b^0 \quad (72)$$

and

$$\frac{\partial \Omega^{1,1}}{\partial F_{m,m}^{1,1}} = -2\rho^{1,1} (A^{0\top} A^0 + \rho^{1,1} F^{1,1\top} F^{1,1})^{-2} \bar{F}_{m,m}^{k,i} A^{0\top} b^0. \quad (73)$$

**Node i)** As shown in Fig. 2i), this node is similar with the node in Fig. 2e). The only difference is that it only have two inputs:  $z^{k,i-1}$  and  $u^{k,i-1}$  for  $k = 1$  and  $i \in \{2,3\}$ . The operation in this layer is:

$$\Omega^{k,i} = (A^{0\top} A^0 + \rho^{k,i} F^{k,i\top} F^{k,i})^{-1} (A^{0\top} b^0 + \rho^{k,i} F^{k,i\top} (z^{k,i-1} - u^{k,i-1})) \quad (74)$$

where

$$A^0 = \frac{1}{2} \bar{\Phi} \quad (75)$$

$$b^0 = \vec{1}. \quad (76)$$

The gradients of output in this layer w.r.t. parameters are:

$$\frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} = -Q^2 F^{k,i\top} F^{k,i} (A^{0\top} b^0 + \rho^{k,i} F^{k,i\top} (z^{k,i-1} - u^{k,i-1})) + Q F^{k,i\top} (z^{k,i-1} - u^{k,i-1}) \quad (77)$$

and

$$\frac{\partial \Omega^{k,i}}{\partial F_{m,m}^{k,i}} = -2\rho^{k,i} Q^2 \bar{F}_{m,m}^{k,i} (A^{0\top} b^0 + \rho^{k,i} F^{k,i\top} (z^{k,i-1} - u^{k,i-1})) + \rho^{k,i} Q \bar{F}_{m,m}^{k,i} (z^{k,i-1} - u^{k,i-1}). \quad (78)$$

The output in this layer w.r.t. its inputs as follows.

$$\frac{\partial \Omega^{k,i}}{\partial z^{k,i-1}} = \rho^{k,i} Q F^{k,i\top} \quad (79)$$

$$\frac{\partial \Omega^{k,i}}{\partial u^{k,i-1}} = -\rho^{k,i} Q F^{k,i\top}. \quad (80)$$

**Node j)** As shown in Fig. 2j), this node is similar with the node in Fig. 2i). The only difference is that its output is the input to compute  $\beta_k$ .

$$\frac{\partial E}{\partial \Omega^{k,i}} = \frac{\partial E}{\partial \beta_k} \frac{\partial \beta_k}{\partial \Omega^{k,i}} \quad (81)$$

for  $k = 1$  and  $i = 4$ .

**Node k):** As shown in Fig. 2k), this node has two sets of inputs:  $u^{k,i-1}$  and  $C^{k,i}$ . Its output  $z^{k,i}$  is the input to compute  $u^{k,i}$  and  $\Omega^{k,i+1}$  for  $k \in \{1,2,3\}$  and  $i \in \{2,3\}$ . Parameter of this layer is  $\rho^{k,i}$ . The operation in this layer is

$$z^{k,i} = S_{\lambda/\rho^{k,i}}(C^{k,i} + u^{k,i-1}), \quad (82)$$

where  $\lambda$  is 1, whose  $n$ th element is

$$z_m^{k,i} = \begin{cases} C_m^{k,i} + u_m^{k,i-1} - 1/\rho^{k,i}, & C_m^{k,i} + u_m^{k,i-1} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i} + u_m^{k,i-1}| \leq 1/\rho^{k,i} \\ C_m^{k,i} + u_m^{k,i-1} + 1/\rho^{k,i}, & C_m^{k,i} + u_m^{k,i-1} < -1/\rho^{k,i}. \end{cases} \quad (83)$$

The gradients of the loss w.r.t. parameters are

$$\frac{\partial E}{\partial \rho^{k,i}} = \frac{\partial E}{\partial z^{k,i}} \frac{\partial z^{k,i}}{\partial \rho^{k,i}}, \quad (84)$$

where

$$\frac{\partial E}{\partial z^{k,i}} = \frac{\partial E}{\partial u^{k,i}} \frac{\partial u^{k,i}}{\partial z^{k,i}} + \frac{\partial E}{\partial \Omega^{k,i+1}} \frac{\partial \Omega^{k,i+1}}{\partial z^{k,i}}. \quad (85)$$

We compute the gradients of the output in this layer w.r.t. parameter as

$$\frac{\partial z_m^{k,i}}{\partial \rho^{k,i}} = \begin{cases} \rho^{k,i-2}, & C_m^{k,i} + u_m^{k,i-1} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i} + u_m^{k,i-1}| \leq 1/\rho^{k,i} \\ -\rho^{k,i-2}, & C_m^{k,i} + u_m^{k,i-1} < -1/\rho^{k,i} \end{cases} \quad (86)$$

and the gradients of the output in this layer w.r.t. the input as

$$\frac{\partial z_m^{k,i}}{\partial C_m^{k,i}} = \frac{\partial z_m^{k,i}}{\partial u_m^{k,i-1}} = \begin{cases} 1, & C_m^{k,i} + u_m^{k,i-1} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i} + u_m^{k,i-1}| \leq 1/\rho^{k,i} \\ 1, & C_m^{k,i} + u_m^{k,i-1} < -1/\rho^{k,i}. \end{cases} \quad (87)$$

**Node l):** As shown in Fig. 2l), this node is similar with the node in Fig. 2k). The only difference is that it only has one input:  $C^{k,i}$  for  $k \in \{1,2,3\}$  and  $i = 1$ . Parameter of this layer is  $\rho^{k,i}$ . The operation in this layer is

$$z^{k,i} = S_{\lambda/\rho^{k,i}}(C^{k,i}), \quad (88)$$

where  $\lambda$  is 1, whose  $n$ th element is

$$z_m^{k,i} = \begin{cases} C_m^{k,i} - 1/\rho^{k,i}, & C_m^{k,i} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i}| \leq 1/\rho^{k,i} \\ C_m^{k,i} + 1/\rho^{k,i}, & C_m^{k,i} < -1/\rho^{k,i}, \end{cases} \quad (89)$$

The gradients of the loss w.r.t. parameters are

$$\frac{\partial E}{\partial \rho^{k,i}} = \frac{\partial E}{\partial z^{k,i}} \frac{\partial z^{k,i}}{\partial \rho^{k,i}}, \quad (90)$$

where

$$\frac{\partial E}{\partial \mathbf{z}^{k,i}} = \frac{\partial E}{\partial \mathbf{u}^{k,i}} \frac{\partial \mathbf{u}^{k,i}}{\partial \mathbf{z}^{k,i}} + \frac{\partial E}{\partial \boldsymbol{\Omega}^{k,i+1}} \frac{\partial \boldsymbol{\Omega}^{k,i+1}}{\partial \mathbf{z}^{k,i}}. \quad (91)$$

We compute the gradients of the output in this layer w.r.t. parameter as

$$\frac{\partial \mathbf{z}_m^{k,i}}{\partial \rho^{k,i}} = \begin{cases} \rho^{k,i-2}, & C_m^{k,i} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i}| \leq 1/\rho^{k,i} \\ -\rho^{k,i-2}, & C_m^{k,i} < -1/\rho^{k,i} \end{cases} \quad (92)$$

and the gradients of the output in this layer w.r.t. the input as

$$\frac{\partial \mathbf{z}_m^{k,i}}{\partial C_m^{k,i}} = \begin{cases} 1, & C_m^{k,i} > 1/\rho^{k,i} \\ 0, & |C_m^{k,i}| \leq 1/\rho^{k,i} \\ 1, & C_m^{k,i} < -1/\rho^{k,i}. \end{cases} \quad (93)$$

**Node m)** As shown in Fig. 2m), this node has one input:  $\Omega^{k,i}$  for  $k \in \{1,2,3\}$  and  $i \in \{1,2,3\}$ . Its output  $C^{k,i}$  is used to compute  $\mathbf{u}^{k,i}$  and  $\mathbf{z}^{k,i}$ . Parameter of this layer is  $F^{k,i}$ . The gradients of loss w.r.t the parameter can be calculated as

$$\frac{\partial E}{\partial F^{k,i}} = \frac{\partial E}{\partial C^{k,i}} \frac{\partial C^{k,i}}{\partial F^{k,i}}, \quad (94)$$

where

$$\frac{\partial E}{\partial C^{k,i}} = \frac{\partial E}{\partial \mathbf{u}^{k,i}} \frac{\partial \mathbf{u}^{k,i}}{\partial C^{k,i}} + \frac{\partial E}{\partial \mathbf{z}^{k,i}} \frac{\partial \mathbf{z}^{k,i}}{\partial C^{k,i}}. \quad (95)$$

The operation in this layer is:

$$C^{k,i} = F^{k,i} \Omega^{k,i}. \quad (96)$$

We compute the gradients of the output in this layer w.r.t parameters

$$\frac{\partial C^{k,i}}{\partial F_{m,m}^{k,i}} = \bar{I}_{m,m}^{k,i} \Omega^{k,i}. \quad (97)$$

The output in this layer w.r.t its inputs are

$$\frac{\partial C^{k,i}}{\partial \Omega^{k,i}} = F^{k,i}. \quad (98)$$

**Node n):** As shown in Fig. 2n), this node has two sets of inputs:  $\beta^{k-1}$  and  $\Omega^{k,i_{max}}$ . Its output  $\beta^k$  is the input to compute  $\Lambda^k$ ,  $\mathbf{v}^k$  and  $\beta^{k+1}$  for  $k = 2$ . Parameters of this layer is  $t^k$ . The operation in this layer is

$$\beta^k = (1 - t^k) \beta^{k-1} + t^k \Omega^{k,i_{max}}. \quad (99)$$

The gradients of loss w.r.t. the parameter can be calculated as

$$\frac{\partial E}{\partial t^k} = \frac{\partial E}{\partial \beta^k} \frac{\partial \beta^k}{\partial t^k}, \quad (100)$$

where

$$\frac{\partial E}{\partial \beta^k} = \frac{\partial E}{\partial \Lambda^k} \frac{\partial \Lambda^k}{\partial \beta^k} + \frac{\partial E}{\partial \mathbf{v}^k} \frac{\partial \mathbf{v}^k}{\partial \beta^k} + \frac{\partial E}{\partial \beta^{k+1}} \frac{\partial \beta^{k+1}}{\partial \beta^k}. \quad (101)$$

We compute the gradients of the output in this layer w.r.t. parameter as

$$\frac{\partial \beta^k}{\partial t^k} = -\beta^{k-1} + \Omega^{k,i_{max}} \quad (102)$$

and the gradients of the output in this layer w.r.t. the input as

$$\frac{\partial \beta^k}{\partial \beta^{k-1}} = (1 - t^k) \mathbf{I}_{M,M} \quad (103)$$

$$\frac{\partial \beta^k}{\partial \Omega^{k,i_{max}}} = t^k \mathbf{I}_{M,M}. \quad (104)$$

**Node o):** As shown in Fig. 2o), this node is similar with the node in Fig. 2n). The only difference is that it only has one input:  $\Omega^{k,i_{max}}$  for  $k = 1$ . Parameters of this layer is  $t^k$ . The operation in this layer is

$$\beta^k = t^k \Omega^{k,i_{max}}. \quad (105)$$

The gradients of the output in this layer w.r.t. parameter as

$$\frac{\partial \beta^k}{\partial t^k} = \Omega^{k,i_{max}} \quad (106)$$

and the gradients of the output in this layer w.r.t. the input as

$$\frac{\partial \beta^k}{\partial \Omega^{k,i_{max}}} = t^k \mathbf{I}_{M,M}. \quad (107)$$

**Node p)** As shown in Fig. 2p), this node is similar with the node in Fig. 2n). The only difference is that its output is not the input of any other nodes for  $k = 3$ . That is  $\frac{\partial E}{\partial \beta^k}$  can be calculated directly:

$$\begin{aligned} \frac{\partial E}{\partial \beta^k} &= \frac{\partial \sum_{n=1}^N \log [1 + \exp \{-y_n \beta^\top \phi(x_n)\}]}{\partial \beta^k} \\ &= - \sum_{n=1}^N y_n \frac{\exp \{-y_n \beta^\top \phi(x_n)\}}{1 + \exp \{-y_n \beta^\top \phi(x_n)\}} \phi(x_n)^\top. \end{aligned} \quad (108)$$

**Node q):** As shown in Fig. 2q), this node has three sets of inputs:  $\mathbf{u}^{k,i-1}$ ,  $C^{k,i}$  and  $\mathbf{z}^{k,i}$ . Its output  $\mathbf{u}^{k,i}$  is the input to compute  $\mathbf{u}^{k,i+1}$ ,  $\mathbf{z}^{k,i+1}$  and  $\Omega^{k,i+1}$  for  $k \in \{1,2,3\}$  and  $i = 2$ . This layer has no parameters. The loss w.r.t. output is

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{u}^{k,i}} &= \frac{\partial E}{\partial \mathbf{u}^{k,i+1}} \frac{\partial \mathbf{u}^{k,i+1}}{\partial \mathbf{u}^{k,i}} \\ &+ \frac{\partial E}{\partial \mathbf{z}^{k,i+1}} \frac{\partial \mathbf{z}^{k,i+1}}{\partial \mathbf{u}^{k,i}} \\ &+ \frac{\partial E}{\partial \Omega^{k,i+1}} \frac{\partial \Omega^{k,i+1}}{\partial \mathbf{u}^{k,i}}. \end{aligned} \quad (109)$$

The operation in this layer is:

$$\mathbf{u}^{k,i} = \mathbf{u}^{k,i-1} + C^{k,i} - \mathbf{z}^{k,i}. \quad (110)$$

The output in this layer w.r.t its inputs are

$$\frac{\partial \mathbf{u}^{k,i}}{\partial \mathbf{u}^{k,i-1}} = \mathbf{I}_{M,M} \quad (111)$$

$$\frac{\partial \mathbf{u}^{k,i}}{\partial C^{k,i}} = \mathbf{I}_{M,M} \quad (112)$$

$$\frac{\partial \mathbf{u}^{k,i}}{\partial \mathbf{z}^{k,i}} = -\mathbf{I}_{M,M}, \quad (113)$$

where  $M$  is the number of elements in  $\beta_k$  and  $\mathbf{I}_{M,M}$  is the identity matrix with the size of  $M \times M$ .

**Node r)** As shown in Fig. 2r), this node is similar with the node in Fig. 2q). The only difference is that its output  $\mathbf{u}^{k,i}$  is the input to compute  $\Omega^{k,i+1}$  for  $k \in \{1,2,3\}$  and  $i = 3$ . The loss w.r.t. output is

$$\frac{\partial E}{\partial \mathbf{u}^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i+1}} \frac{\partial \Omega^{k,i+1}}{\partial \mathbf{u}^{k,i}}. \quad (114)$$

**Node s)** As shown in Fig. 2t), this node is similar with the node in Fig. 2q). The only difference is that it has only two inputs:  $\mathbf{C}^{k,i}$  and  $\mathbf{z}^{k,i}$  for  $k \in 1,2,3$  and  $i = 1$ . The operation in this layer is:

$$\mathbf{u}^{k,i} = \mathbf{C}^{k,i} - \mathbf{z}^{k,i}. \quad (115)$$

The output in this layer w.r.t its inputs are

$$\frac{\partial \mathbf{u}^{k,i}}{\partial \mathbf{C}^{k,i}} = \mathbf{I}_{M,M} \quad (116)$$

$$\frac{\partial \mathbf{u}^{k,i}}{\partial \mathbf{z}^{k,i}} = -\mathbf{I}_{M,M}. \quad (117)$$

**Gradients of parameters shared among layers):** As we have discussed in node e), parameters may share among different layers. In the above calculations, we only calculate the gradients of parameters dependent on the output of this node. To get a complete expression of gradients, we need to sum them together. Thus, we have

$$\frac{\partial E}{\partial \rho^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} + \frac{\partial E}{\partial \mathbf{z}^{k,i}} \frac{\partial \mathbf{z}^{k,i}}{\partial \rho^{k,i}} \quad (118)$$

$$\frac{\partial E}{\partial \mathbf{F}^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \mathbf{F}^{k,i}} + \frac{\partial E}{\partial \mathbf{C}^{k,i}} \frac{\partial \mathbf{C}^{k,i}}{\partial \mathbf{F}^{k,i}} \quad (119)$$

for  $i \in \{1,2,3\}$ , and

$$\frac{\partial E}{\partial \rho^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \rho^{k,i}} \quad (120)$$

$$\frac{\partial E}{\partial \mathbf{F}^{k,i}} = \frac{\partial E}{\partial \Omega^{k,i}} \frac{\partial \Omega^{k,i}}{\partial \mathbf{F}^{k,i}} \quad (121)$$

for  $i = 4$ .

## REFERENCES

[1] J. R. Hershey, J. L. Roux, and F. Weninger, "Deep unfolding: Model-based inspiration of novel deep architectures," *arXiv preprint arXiv:1409.2574*, 2014.

[2] D. D. Lee and H. S. Seung, "Algorithms for non-negative matrix factorization," *Advances in Neural Information Processing Systems*, 2001, pp. 556–562.

[3] S. Wisdom, J. Hershey, J. Le Roux, and S. Watanabe, "Deep unfolding for multichannel source separation," in *Proc. 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 121–125.

[4] C.-H. Lee and J.-T. Chien, "Deep unfolding inference for supervised topic model," in *Proc. 2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, IEEE, 2016, pp. 2279–2283.

[5] R. Vemulapalli, O. Tuzel, and M.-Y. Liu, "Deep gaussian conditional random field network: A model-based deep network for discriminative denoising," in *Proc. the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4801–4809.

[6] H. He, B. Xin, and D. Wipf, "From sparse bayesian learning to deep recurrent nets," *NIPS*, 2017.

[7] D. Mahapatra, S. Mukherjee, and C. S. Seelamantula, "Deep sparse coding using optimized linear expansion of thresholds," *arXiv preprint arXiv:1705.07290*, 2017.

[8] J. Sun, H. Li, Z. Xu *et al.*, "Deep admm-net for compressive sensing mri," *Advances in Neural Information Processing Systems*, pp. 10–18, 2016.

[9] M. E. Tipping, "Sparse bayesian learning and the relevance vector machine," *Journal of machine learning research*, vol. 1, no. Jun, pp. 211–244, 2001.

[10] W. Pan, "Bayesian learning for nonlinear system identification," Ph.D. dissertation, Imperial College London, 2015.

[11] Y. Li, C. Campbell, and M. Tipping, "Bayesian automatic relevance determination algorithms for classifying gene expression data," *Bioinformatics*, vol. 18, no. 10, pp. 1332–1339, 2002.

[12] B. Krishnapuram, L. Carin, and A. J. Hartemink, "Joint classifier and feature optimization for comprehensive cancer diagnosis using gene expression data," *Journal of Computational Biology*, vol. 11, no. 2-3, pp. 227–242, 2004.

[13] G. C. Cawley and N. L. Talbot, "Gene selection in cancer classification using sparse logistic regression with bayesian regularization," *Bioinformatics*, vol. 22, no. 19, pp. 2348–2355, 2006.

[14] X. Yang, W. Pan, and Y. Guo, "Sparse bayesian classification and feature selection for biological expression data with high correlations," *PLOS One*, vol. 12, no. 12, 2017.

[15] A. L. Yuille and A. Rangarajan, "The concave-convex procedure," *Neural Computation*, vol. 15, no. 4, pp. 915–936, 2003.

[16] S.-I. Lee, H. Lee, P. Abbeel, and A. Y. Ng, "Efficient l1 regularized logistic regression," in *AAAI*, vol. 6, 2006, pp. 401–408.

[17] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proc. COMPSTAT'2010*. Springer, 2010, pp. 177–186.

[18] T. P. Minka. (2003). A comparison of numerical optimizers for logistic regression. [Online]. Available: <http://research.microsoft.com/~minka/papers/logreg>

[19] T. Lipp and S. Boyd, "Variations and extension of the convex-concave procedure," *Optimization and Engineering*, vol. 17, no. 2, pp. 263–287, 2016.

[20] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and TrendsR in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.

[21] S. Yan, X. Yang, C. Wu, Z. Zheng, and Y. Guo, "Balancing the stability and predictive performance for multivariate voxel selection in fmri study," in *Proc. International Conference on Brain Informatics and Health*. Springer, 2014, pp. 90–99.

[22] P. Kirk, A. Witkover, C. R. Bangham, S. Richardson, A. M. Lewin, and M. P. Stumpf, "Balancing the robustness and predictive performance of biomarkers," *Journal of Computational Biology*, vol. 20, no. 12, pp. 979–989, 2013.

[23] K. Fernandes, J. S. Cardoso, and J. Fernandes, "Transfer learning with partial observability applied to cervical cancer screening," in *Proc. Iberian Conference on Pattern Recognition and Image Analysis*. Springer, 2017, pp. 243–250.

[24] S. L. Pomeroy, P. Tamayo *et al.*, "Prediction of central nervous system embryonal tumour outcome based on gene expression," *Nature*, vol. 415, no. 6870, pp. 436–442, 2002.

[25] Z. Wu, R. A. Irizarry, R. Gentleman, F. Martinez-Murillo, and F. Spencer, "A model-based background adjustment for oligonucleotide expression arrays," *Journal of the American statistical Association*, vol. 99, no. 468, pp. 909–917, 2004.



**Xian Yang** received her PhD degree from the Department of Computing, Imperial College London in 2016. She has been working as a research assistant at Imperial College London since 2012 and became a research associate in 2016. Her research interests include machine learning, neuroimaging, bioinformatics, system biology, data mining, statistics and health informatics. She has taken part in various projects, such as the UBIOPRED project on the severe asthma study, the iHealth project on clinical pathway management, and the eTRIKS project on knowledge management for precision medicine research.



**Yike Guo** is a professor of computing science in the Department of Computing at Imperial College London. He is the founding director of the Data Science Institute at Imperial College, as well as leading the Discovery Science Group in the department. He has been working on technology and platforms for scientific data analysis, where his research focuses on knowledge discovery, data mining and large-scale data management. Professor Guo also holds the position of CTO of the transSMART Foundation, a global open source community using and developing data sharing and analytics technology for translational medicine, and the position of CIO of IDBS, a world's leading company in developing innovative data management and analytics solutions.