

A Mixture of Views Network with Applications to the Classification of Breast Microcalcifications

Yaniv Shachor¹ Hayit Greenspan² Jacob Goldberger¹

¹ Faculty of Engineering, Bar-Ilan University, Ramat-Gan, Israel

² Department of Biomedical Engineering, Tel Aviv University, Tel Aviv, Israel.

Abstract. In this paper we examine data fusion methods for multi-view data classification. We present a decision concept which explicitly takes into account the input multi-view structure, where for each case there is a different subset of relevant views. This data fusion concept, which we dub Mixture of Views, is implemented by a special purpose neural network architecture. It is demonstrated on the task of classifying breast microcalcifications as benign or malignant based on CC and MLO mammography views. The single view decisions are combined by a data-driven decision, according to the relevance of each view in a given case, into a global decision. The method is evaluated on a large multi-view dataset extracted from the standardized digital database for screening mammography (DDSM). The experimental results show that our method outperforms previously suggested fusion methods.

1 Introduction

Over the years, research has been actively seeking the best ways to integrate data from multiple sources. Data fusion can be performed at the raw data, feature level and decision level. Feature-level fusion involves concatenating the features extracted from each view as an input to a decision system. Decision level fusion is based on averaging the decisions obtained from each sensor independently.

In the medical field, multi view data fusion has been researched for number of applications. Prasoon et al. [6] combined features from orthogonal patches in order to segment knee cartilage. Setio et al. [8] used decision level fusion on orthogonal patches in a pulmonary nodule detection CAD. The current study deals with multi-view mammography image information fusion. In this task, information from images acquired from different angles is aggregated for automatic classification of breast microcalcifications (MC) as benign or malignant. Recent studies (e.g. [7] [10]) have confirmed the superior performance of a multi-view CADx system over its single-view counterpart. Bekker et al. [1] proposed averaging the two view-level decisions. There is, however, a structure specific to this problem that is not explicitly modeled by equal weight averaging or even by a weighted average with fixed weights. Screening mammography typically involves taking two views of the breast, from above and from an oblique, since a 3D pathology indication is not always clearly observed in a single 2D image. Thus an expert radiologist's diagnosis of malignance is based on evidence that is not

necessarily clearly seen in both views.

In this study we introduce a neural network (NN) architecture that explicitly takes into account the multiple-view structure of the problem by integrating view level decision making. We present an automatic data-driven strategy that finds which view conveys more relevant information for clinical decision making. Instead of simple averaging of the view-level decisions, we train a ‘gating’ network that decides the best way to average the view-level decisions for each case. Our method is related to the well-known Mixture of Experts (MoE) model [4]. Unlike classical MoE which divides the task among a set of experts in an unsupervised way, in MoV each expert is associated with a sensor (or a view). We describe a neural network architecture that implements the MoV concept.

To evaluate the method we use the labeled multi-view mammogram dataset DDSM [3]. It contains MCs location in both views marked by experts and we also have the biopsy results, showing whether the abnormalities were benign or malignant. Experiments were performed on pairs of CC+MLO views extracted from the DDSM dataset. Fig. 1 shows benign and malignant examples from the DDSM dataset. The results when applying the MOV to the DDSM dataset show that using our approach outperforms previously suggested fusion methods.

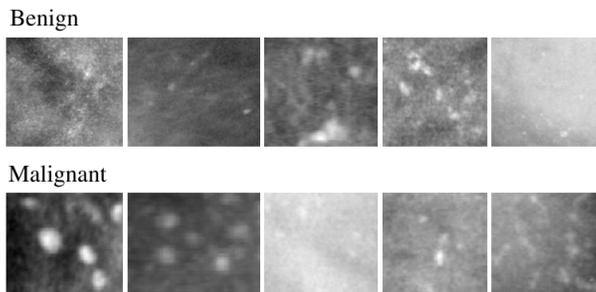


Fig. 1: Examples of benign and malignant MC clusters (from the DDSM dataset). ROI sizes range from $2mm \times 2mm$ to $6mm \times 6mm$.

2 The Mixture of Views Network architecture

In this study we deal with a classification task where the features are obtained from multiple sensors. Our goal is to construct a neural network architecture which is aware of this input structure and takes advantage of it to improve classification performance. In this section we first describe the probabilistic framework we use to model the multi-view data fusion classification and then derive a training algorithm that simultaneously finds the parameters of each view-based classifier and the parameters of a gating network that decides which view is more relevant for a given input feature set. Hereafter we use the terms sensor and view interchangeably.

Assume a feature vector x is a concatenation of m components $x = (x^1, \dots, x^m)$ such that each component is obtained from a different sensor. Each sensor can provide a different number of features. The standard way of classifying an object based on features from multiple views is to concatenate the view-level inputs and using the concatenated vector as an input to a standard neural-network classifier.

A model based on concatenation, however, does not take into account the structure of the system, as a fusion of several views. Each view has enough information to make its own reasonable prediction and, in principle, each component can be used alone for the classification task. In addition not all the views are equally relevant in each case. We propose a model that computes each view-level prediction separately and then combines them to form the final estimation. In other words, in our approach we concatenate the view-level decisions instead of the view-level features. We also analyze the features provided by different views and combine the predictions dynamically, according to the relevance of the view in each case.

Assume we are given a k -class classification problem with labels denoted by $1, \dots, k$ and input x composed of m view-level components denoted by x^1, \dots, x^m . For the classification task we use a neural network that combines the view-level decisions. Let $\theta = \{\theta_g, \theta_1, \theta_2, \dots, \theta_m\}$ be the model parameters, where θ_i are the parameters of the network that performs a classification based solely on the component x^i and θ_g is the parameter-set of a gating network that computes a data-driven distribution over all the views. The final decision is derived by computing a weighted average of the view-level decisions where the weights are provided by the gating networks.

The network has three main components: a set of neural-networks where each one renders a decision based solely on a single sensor, a gate neural-network that defines those sensors whose individual sensor-based opinions are trustworthy, and a weighted sum of experts, where the weights are the gating network outputs.

In the MoV model the probability of input x being labeled as $c \in \{1, \dots, k\}$ is:

$$p(y = c|x; \theta) = \sum_{i=1}^m p(i|x; \theta_g) p(y = c|x^i; \theta_i). \quad (1)$$

We can view the model as a two-step process that produces a decision y given an input feature set x . We first use the gating function to select the view that conveys the relevant information for decision making and then apply the corresponding network to determine the output label.

As stated above, the MoV model can be viewed as an instance of Mixture of Experts (MoE) modeling. The MoE approach was introduced more than twenty years ago [4], and combines the decisions of several experts, each of which specializes in a different part of the input space. The MoE model is based on the “divide and conquer” paradigm, which solves complex problems by dividing them into simpler ones and combine their solutions to solve the original task. The model allows the individual experts to specialize on smaller parts of a larger problem, and it uses soft partitions of the data implemented by the gate. In the general setup of MoE all experts are exposed to all the features and the goal is to

cluster the feature space and associate each cluster with an expert classifier in an unsupervised manner. By contrast, in the MoV model there is a predefined partitioning of the set of features according to the sensors used to measure them and each expert is specialized in making decisions based solely on the features of the corresponding sensor.

We next describe the training procedure. Assume we are given n feature vectors x_1, \dots, x_n with corresponding labels $y_1, \dots, y_n \in \{1, \dots, k\}$. Each input vector x_i is composed of m components x_i^1, \dots, x_i^m collected from the m sensors. The log-likelihood function of the model parameters is:

$$L(\theta) = \sum_t \log p(y_t|x_t; \theta) = \sum_t \log \left(\sum_i p(i|x_t; \theta_g) p(y_t|x_t^i; \theta_i) \right). \quad (2)$$

To find the network parameters we can maximize the likelihood function using the standard back-propagation algorithm. It can be easily verified that the back-propagation equation for the parameter set of the i -th expert is:

$$\frac{\partial L}{\partial \theta_i} = \sum_{t=1}^n w_{ti} \cdot \frac{\partial}{\partial \theta_i} \log p(y_t|x_t^i; \theta_i) \quad (3)$$

such that w_{ti} is the posterior distribution of the gating decision:

$$w_{ti} = p(i|x_t, y_t; \theta) = \frac{p(y_t|x_t^i; \theta_i) p(i|x_t; \theta_g)}{p(y_t|x_t; \theta)}. \quad (4)$$

In a similar way, the back-propagation equation for the parameter set of the gating network is:

$$\frac{\partial L}{\partial \theta_g} = \sum_{t=1}^n \sum_{i=1}^m w_{ti} \cdot \frac{\partial}{\partial \theta_g} \log p(i|x_t; \theta_g). \quad (5)$$

The likelihood score (2) is focused on the performance of the compound network. It does not, however, explicitly encourage each view-level network to obtain the optimal decision based on features of the corresponding view. One strategy to overcome this issue is to train initially each view-level network separately and then to use the learned parameters as initial values for the compound network. We have found that this approach, which is based on injecting information to the system via parameter initialization, does not work well since the network tends to forget its initialization after a few training iterations. Rather, we use a modified likelihood score:

$$L(\theta) + \lambda \sum_i L_i(\theta_i) \quad (6)$$

such that $L(\theta)$ is the usual likelihood score (2) and $L_i(\theta_i) = \sum_t \log p(y_t|x_t^i; \theta_i)$ is the likelihood score where we only use the network corresponding to the i -th view for classification. The parameter λ controls the relative importance of the view-level and integrated decisions and can be tuned using cross-validation.

3 Multi-view Classification of Breast Microcalcifications

In this section we demonstrate the MoV method in the task of classifying breast microcalcifications as benign or malignant, based on two mammography views. A screening mammographic examination usually consists of four images, corresponding to each breast scanned in two views: the mediolateral oblique (MLO) view and the craniocaudal (CC) view. The MLO projection is taken in a 45° angle and shows part of the pectoral muscle. The CC projection is a top-down view of the breast. Both views are included in the diagnostic procedure. When reading mammograms, radiologists judge whether or not a malignant lesion is present by examining both views and breasts. In an expert diagnosis procedure, the expert looks at each of the views separately, and delivers one final assessment.

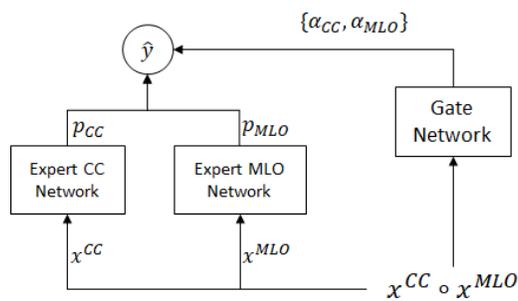


Fig. 2: MoV network architecture for multi-view classification of breast microcalcifications.

We next apply the MoV method presented above for classification of breast microcalcifications as benign or malignant. In this task, let $\{x^{cc}, x^{mlo}\}$ denote the extracted features from the CC and MLO views, respectively. The MoV network for this task is:

$$p(c|x; \theta) = \sum_{i \in \{cc, mlo\}} p(i|x^{cc} \circ x^{mlo}; \theta_g) p(c|x^i; \theta_i) \quad (7)$$

such that c is either benign or malignant. The network is illustrated in Fig. 2. In the next section we show empirically that the MoV network yields better classification results than networks that are based on a single view and better results than other strategies that combine information from the two views.

4 Experimental Setup

Dataset and features. The empirical evaluation was based on the DDSM dataset [3] which provides a high number of annotated mammograms with a

biopsy-proven diagnosis. We extracted ROIs that contained clusters of MCs for which a proven pathology had been found. In real-life diagnosis procedures, radiologists pay attention to the density of the breast, which can be a good predictor of a woman’s breast cancer risk. We assumed that categorizing the mammograms based on their density was necessary to compare the different features and classifiers in a more objective way. We separated the mammograms into two different tissue-density categories and studied them individually: fatty tissues (ratings 1 and 2), and dense tissues (ratings 3 and 4). We considered classifying the fatty and dense breast densities as two separate tasks. The density type for each case was a parameter supplied by an expert as part of the DDSM. We chose patients in the DDSM dataset that had both CC and MLO views to test our model. Our dataset was comprised of 1410 clusters (705 of CC, and 705 of MLO), of which 372 were benign and 333 were malignant.

Feature vectors x_{cc} and x_{mlo} were extracted from the CC and MLO views, respectively. Following [1], texture features were extracted from the Curvelet coefficients at intermediate scales (in our study, two scales), and included the four features mentioned in [9] for each scale, with three additional features: entropy, skewness and kurtosis. Overall, each extracted ROI was represented by 14 features.

Training procedure. We used a 10-fold cross validation setup. In this setup, there is complete isolation of the test set from the train set. Each fold was only used for testing and never for training. We thus ensured that no bias was introduced. In addition, 10% from each training set was used as a validation set in order to optimize the model hyper-parameters, according to the mean results on the validation sets. Using the features described in the previous section, the size of the input feature set is 28 (14 features for each view). The features of each view were inserted into the expert NN. In addition, all the features were inserted into the gate NN, to receive the weights for the experts’ predictions. The expert NN has 2 hidden layers comprised of 24 neurons each. The gate NN has 2 hidden layers, comprised of 3 neurons each. We used ReLU non-linear activations between the layers and dropout with parameter 0.5. We used Eq. (6) as a target function for optimization where $\lambda = 1$ gave the best results. The objective function was maximized using the Adam optimizer [5]. The network trained over 500 epochs with reduction of the learning rate on the loss plateau and early stopping.

Table 1: Classification results (benign vs. malignant) for fatty (left) and dense (right) breast tissues.

Method	Accuracy	F-measure	AUC
CC	0.704	0.630	0.695
MLO	0.698	0.583	0.682
Avg [1] [8]	0.709	0.620	0.699
concat [6]	0.718	0.624	0.704
MoV	0.718	0.632	0.708

Method	Accuracy	F-measure	AUC
CC	0.643	0.580	0.639
MLO	0.629	0.520	0.622
Avg [1] [8]	0.646	0.573	0.641
Concat [6]	0.642	0.568	0.637
MoV	0.665	0.596	0.661

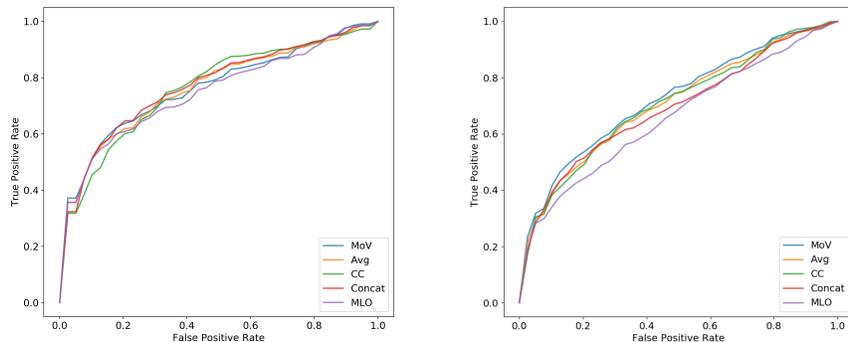


Fig. 3: ROC curves for the fatty (left) and dense (right) tissues.

Experimental evaluation. We compared the MoV method to several models: a network based on the single-view (CC or MLO) features as input, and a network that averages the decisions based on CC and MLO [1] [8]. This model can be viewed as a degenerated version of MoV where the data-driven gating network is replaced by a simple averaging. We denote this model as Avg. We also implemented a network that concatenates CC and MLO features as input, as in [6], denoted as Concat.

We optimized each model’s hyper-parameters using a validation set, to conduct a fair comparison. We used Receiver Operator Characteristic (ROC) curves with accuracy, Area Under Curve (AUC) and F-measure metrics to evaluate the algorithms’ performance. Due to the differences in nature between the two types of breast tissue (Fatty/Dense), a different network architecture was trained for each type, and therefore each breast tissue category was evaluated separately. Table 1 shows the classification results for the two types of breast tissue (Fatty/Dense), and Fig. 3 shows the corresponding ROC curves. The results are the 10-fold test set classification average, computed over several experiments. As can be seen from the tables, using the MoV method yielded the best results over the dense samples in all metrics and over the fatty samples, had the best results on the F-measure and AUC and had the best accuracy on a par with the Concat network.

We performed DeLong test [2] comparing the MoV model paired with each of the baseline models. The input to the DeLong test consisted of predictions from the 10-fold cross validation with the corresponding labels. The DeLong test examined the null hypothesis that both methods have the same AUC. On the dense data, all the hypotheses were successfully rejected with p-value < 0.05 . On the fatty data we have not seen significance in the results. We assume it is because the fatty data has half of the size of the dense dataset.

To conclude, in this study we addressed the problem of fusing several data sources for a classification task. We proposed a network architecture that is explicitly aware that the input data are provided by multiple sensors. We demonstrated the performance of the MoV method on the classification of breast microcalcifications into benign and malignant given multi-view mammograms. We showed that the MoV architecture yields improved performance. In the algorithm presentation we focused on the a two-view mammography, namely CC and MLO. Our method, however, can be easily extended to mammography with more than two views. In addition, the ROI features could be extracted automatically using transfer learning methods. In the future we plan to investigate the applicability of the method to other medical imaging tasks with multiple views and/or multiple scans.

References

1. Bekker, A.J., Shalton, M., Greenspan, H., Goldberger, J.: Multi-view probabilistic classification of breast microcalcifications. *IEEE Trans. Medical Imaging* 35:2, 645–653 (2016)
2. DeLong, E.R., DeLong, D.M., Clarke-Pearson, D.L.: Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach. *Biometrics* 44(3), 837–845 (1988)
3. Heath, M., Bowyer, K., Kopans, D., Moore, R., Kegelmeyer, W.P.: The digital database for screening mammography. *Proceedings of the Fifth International Workshop on Digital Mammography*, M.J. Yaffe, ed pp. 212–218 (Medical Physics Publishing, 2001)
4. Jacobs, R.A., Jordan, M.I., Nowlan, S.J., Hinton, G.E.: Adaptive mixtures of local experts. *Neural Computation* 3, 79–87 (1991)
5. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *CoRR* abs/1412.6980 (2014)
6. Prasoon, A., Petersen, K., Igel, C., Lauze, F., Dam, E., Nielsen, M.: Deep feature learning for knee cartilage segmentation using a triplanar convolutional neural network. In: *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. pp. 246–253 (2013)
7. Samulski, M., Karssemeijer, N.: Optimizing case-based detection performance in a multi-view CAD system for mammography. *IEEE Trans on Med Imaging* 30, 1001–1009 (2011)
8. Setio, A.A.A., Ciompi, F., Litjens, G., Gerke, P., Jacobs, C., van Riel, S.J., Wille, M.M.W., Naqibullah, M., Snchez, C.I., van Ginneken, B.: Pulmonary nodule detection in CT images: False positive reduction using multi-view convolutional networks. *IEEE Transactions on Medical Imaging* 35(5), 1160–1169 (2016)
9. Shang, Y., Diao, Y., Li, C.: Rotation invariant texture classification algorithm based on curvelet transform and SVM. In: *Machine Learning and Cybernetics, 2008 International Conference on*. vol. 5, pp. 3032–3036 (2008)
10. Velikova, M., Lucas, P., Smulski, M., Karssmeijer, N.: A probabilistic framework for image information fusion with an application to mammographic analysis. *Medical Image Analysis* 16, 865–875 (2012)