

# Regression with Gaussian Mixture Networks

In this document, we provide technical details about the discriminative model  $P^d$  and its optimisation. Recall that  $P^d(\vec{X}_k|\vec{X}_k \neq 0)$  is a Gaussian Mixture probability distribution in  $\mathbb{R}^2$ , where the weight of each Gaussian is predicted by a Neural Network.

We therefore need to learn the parameters of the probability distribution

$$P^d(\vec{X}_k|\vec{X}_k \neq 0) = \sum_{1 \leq m \leq M} f_h(\mathcal{F}_k^c; \theta_h)_m \mathcal{N}(\vec{X}_k - \alpha_m; \sigma_m), \quad (1)$$

namely, the Gaussian parameters  $\alpha_m$  and  $\sigma_m$  for each mode index  $m$ , and the network parameters  $\theta_h$ .

Following Eq.17 of the main paper, we treat each pixel as an independent data-point  $\vec{x}_s$ .  $S = S(Z^0, \dots, Z^D)$  denotes the set of those data-points. In this section, we assume that we have access to a label in  $\mathbb{R}^2$ , for each data-point  $\vec{x}_s$ . We recall that this label is generated by sampling of the generative model, given *ground truth* detections  $Z$ .

The procedure that we use to optimise the following loss derived from Eq.17

$$\mathcal{R}(\theta_h, \alpha, \sigma) = - \sum_{(\vec{x}_s \in S)} \log(P^d(\vec{x}_s|\mathcal{F}_{k_s}^{c_s}, \theta_h, \alpha, \sigma)), \quad (2)$$

follows the same principles as the standard Gaussian Mixture regression model via Expectation-Maximization algorithm [1] and is also closely related to the recent Neural Decision Forests [2], which introduce a Network producing a probability distribution in the form of a mixture of Histograms.

**Updating the Network Parameters** We update the parameters of the network  $\theta_h$  by direct back-propagation and stochastic gradient descent on the objective of Eq. 2.

**Updating the Gaussian Parameters** Let  $\alpha^t$  and  $\sigma^t$  denote the current estimates of the the Gaussian parameters. We derive a closed form update which guarantees that

$$\mathcal{R}(\theta_h, \alpha^{t+1}, \sigma^{t+1}) \leq \mathcal{R}(\theta_h, \alpha^t, \sigma^t). \quad (3)$$

For each data point  $\vec{x}_s$ , let us introduce the distribution over the mixture elements  $m$ ,

$$\xi^t(m|\vec{x}_s, \mathcal{F}_{k_s}^{c_s}, \theta_h, \alpha^t, \sigma^t) = \frac{f_h(\mathcal{F}_{k_s}^{c_s}; \theta_h)_m \mathcal{N}(\vec{x}_s - \alpha_m^t; \sigma_m^t)}{\sum_{1 \leq m' \leq M} f_h(\mathcal{F}_{k_s}^{c_s}; \theta_h)_{m'} \mathcal{N}(\vec{x}_s - \alpha_{m'}^t; \sigma_{m'}^t)} \quad (4)$$

usually called “responsibilities” in the GMM literature.

We then use the standard variational trick with the auxiliary distribution  $\xi^t(m)$  to minimise an upper-bound on  $\mathcal{R}(\theta_h, \alpha, \sigma)$ , with respect to the parameters  $\alpha$  and  $\sigma$ .

$$\begin{aligned}
\mathcal{R}(\theta_h, \alpha, \sigma) &= - \sum_{\vec{x}_s \in S} \log \left( \sum_{1 \leq m \leq M} f_h(\mathcal{F}_{k_s}^{c_s}; \theta_h)_m \mathcal{N}(\vec{x}_s - \alpha_m; \sigma_m) \right) \\
&= - \sum_{\vec{x}_s \in S} \log \left( \sum_{1 \leq m \leq M} \xi^t(m|\vec{x}_s) \frac{f_h(\mathcal{F}_{k_s}^{c_s}; \theta_h)_m \mathcal{N}(\vec{x}_s - \alpha_m; \sigma_m)}{\xi^t(m|\vec{x}_s)} \right) \\
&\leq - \sum_{\vec{x}_s \in S} \sum_{1 \leq m \leq M} \xi^t(m|\vec{x}_s) \log \left( \frac{f_h(\mathcal{F}_{k_s}^{c_s}; \theta_h)_m \mathcal{N}(\vec{x}_s - \alpha_m; \sigma_m)}{\xi^t(m|\vec{x}_s)} \right) \\
&\leq \mathcal{R}(\theta_h, \alpha^t, \sigma^t) - \sum_{\vec{x}_s \in S} \sum_{1 \leq m \leq M} \xi^t(m|\vec{x}_s) \log \left( \frac{\mathcal{N}(\vec{x}_s - \alpha_m; \sigma_m)}{\mathcal{N}(\vec{x}_s - \alpha_m^t; \sigma_m^t)} \right) \quad (5)
\end{aligned}$$

Minimizing Eq. 5 with respect to  $\alpha$  and  $\sigma$  is a convex problem. Assuming that we can find the values achieving the minimum, let us set  $\alpha^{t+1}$  and  $\sigma^{t+1}$  to these. Then, from Eq. 5, we obtain

$$\mathcal{R}(\theta_h, \alpha^{t+1}, \sigma^{t+1}) \leq \mathcal{R}(\theta_h, \alpha^t, \sigma^t),$$

with equality if and only if  $\alpha^{t+1} = \alpha^t$  and  $\sigma^{t+1} = \sigma^t$ .

We therefore need to minimize Eq. 5 with respect to  $\alpha$  and  $\sigma$ , which is equivalent to maximizing

$$\sum_{\vec{x}_s \in S} \sum_{1 \leq m \leq M} \xi^t(m|\vec{x}_s) \log(\mathcal{N}(\vec{x}_s - \alpha_m; \sigma_m)),$$

with respect to the parameters  $\alpha$  and  $\sigma$ . This is done by using the standard optimality conditions for convex problems. We obtain

$$\alpha_m^{t+1} = \frac{\sum_{\vec{x}_s \in S} \xi^t(m|\vec{x}_s) \vec{x}_s}{\sum_{\vec{x}_s \in S} \xi^t(m|\vec{x}_s)}, \quad (6)$$

and,

$$\sigma_m^{t+1} = \frac{\sum_{\vec{x}_s \in S} \xi^t(m|\vec{x}_s) (\vec{x}_s - \alpha_m^{t+1})^2}{\sum_{\vec{x}_s \in S} \xi^t(m|\vec{x}_s)}. \quad (7)$$

**Alternating both** In practice, we alternate one epoch of stochastic gradient descent optimizing the network parameters  $\theta_h$  with one update of the Gaussian parameters of Eqs. 6 and 7. For memory usage reasons, the sums in Eqs. 6 and 7, have to be split into mini-batches. However the update is done after summation over the whole dataset or a very large number of samples.

## References

- [1] C. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006. 1
- [2] P. Kotschieder, M. Fiterau, A. Criminisi, and S. R. Buló. Deep neural decision forests. In *International Conference on Computer Vision*, 2015. 1