

# Gaussian parametrization of deep convolutional network filters

Domen Tabernik

CN-CR Centre, School of Computer Science, University of Birmingham, Birmingham, UK

*Abstract—*

## I. MOTIVATION

Deep neural networks have achieved remarkable progress in last years, outperforming many other methods on different computer vision problems. In particular, convolutional variant of neural network enabled visual object recognition on higher resolution images, ReLU non-linearity reduced learning time, and improved regularization with weight decay and drop-out overcame the strong overfitting issues common with neural networks. Nevertheless, representations being learned by ConvNets are not well understood. Krizhevsky et al. [1] showed first-layer filters resemble various blob detectors and Gabor filters of different orientations. Zeiler and Fergus [2] further showed corner and edge/color conjunctions appear on the second layer, a more complex texture pattern on the third layer and class-specific parts, such as dog's face or bird's legs, on the fourth layer. While Mahendran and Vedaldi [3] found lower layers to be more accurate in reconstruction of the image, they also found higher layers invert back to compositions of parts similar to, but not identical to, the ones in the original image.

With the exception of the first layers, such visualization techniques point to neural networks encoding features as object parts that are composed from simpler parts and partially correlate with semantic parts of an object. This is partially inherent in the network with a deep, multilayer organization allowing complex features to be composed from simpler ones. However, filters that define features are free to learn any representation they find the best for the optimization criteria used, and the most optimal one for object representation seems to capture features as a composition of a partially semantic parts. That holds true at least to certain degree since not all features will correlate well with semantic parts and, as pointed before, first layers capture mostly blob-like and edge features. Such representation also appears quite natural since all objects we observe have a natural hierarchy of compositions. For instance, a visual representation of a cat is composed from a head, a tail, a couple of legs, and a fur. Furthermore, each one of those can be further composed from yet more simple features. A head is composed from ears, whiskers, two eyes, a mouth and a nose. All of those features must not only be present, but appear in specific spatial relation. This represents a clear structure that defines each object.

Despite such a clear structure present in visual domain neural networks do not directly exploit it and do not constrain

to specific structure in the data. The only constraint is the assumption of hierarchical structure enforced through deep, multilayer organization. This makes deep networks ideal for problems where very little can be said about the representation, but in visual domain much can be said about it. One such prior knowledge has already been incorporated by LeCun [4] with weight-sharing enforcing the translation invariance in images. However, ConvNet imposes only hierarchical structure but no spatial structure since spatial relation between sub-features is undefined. Feature's response is depended on having underlying sub-features present in the specific pattern. This is governed by a filter which is learned pixel-by-pixel. The size of the filter defines the how much of the local neighborhood needs to be considered, but specific spatial relation between a feature and its sub-feature is unconstrained since filter is unstructured, and can take any arbitrary form. This can be important in first layer where Gabor filters can be trained without specifying that we are looking for them, but in higher layers it is oblivious of the existing spatial structure. This prevents any kind of semantical interpretation of features as well, since filters are only a bunch of pixel values without any meaning. Despite of freedom in filter learning many filters end up having smooth values with clearly visible spatial structures. Often smoothness of filters is also a requirement to get good results. Furthermore, learning filters directly in pixel-by-pixel manner adds significant amount of optimization parameters, which in return requires a huge amount of data to avoid overfitting and at the same time limits the size of filters that can be used.

In this work we propose to extend the convolutional networks with the assumption of spatial structure present in visual objects. We constrain the structure of filters in convolutional networks to better capture that. This is achieved by parametrizing the filters with Gaussian distributions. Instead of learning individual pixel values in filters we allow Gaussian distributions to govern the filters. This has multiple implications. Firstly, each feature can now be considered as a composition of sub-features where multiple Gaussian components define its characteristics. Instead of defining an occurrence pattern of a sub-feature with an unconstrained filter it can now be defined with multiple Gaussian components. On one hand this still allows to have an arbitrary filter definition when over-parametrizing with many components per filter, but on the other hand a sparse solution would allow to have a compositional interpretation, since Gaussian components can be seen as definition of sub-feature's characteristics. In fact, Gaussian component directly governs strength as well as spatial position and variability of a sub-feature, and different configurations of components represent different filters, such as, blob and

edge detectors on lower layer, or compositions on higher layers. A second implication of Gaussian parametrization is smoothness of filters. Normally this property is an indication of correct learning and parametrization now directly enforces it. This is done through the variance in Gaussian distribution which directly relates to the smoothness of the filter. The last implication is also a reduction in the number of parameters per filter and thus reduction in parameter search-space. A filter of size  $9 \times 9$  can be described by  $3 \times 3$  Gaussian components. Assuming Gaussian components are described with four parameters, weight, mean position in two dimensions and variance, this reduces the number of parameters from 81 to 36, and can be further reduced if sparse solution is required.

In this work we focus on showing how parametrization with Gaussian distribution can be incorporated into convolutional neural network. We further analyze how gradient descent and back-propagation can be applied to learn Gaussian parameters, and if it can be successfully applied to image classification. Furthermore, we show ConvNet with filter parametrization using Gaussian distribution can be on one hand considered as a regular convolutional network, while on the other hand it can be considered as a compositional hierarchical model. That means it can act as a bridge between convolutional networks and compositional hierarchies, and opens an opportunity in future to further improve both of them. For instance, powerful optimization with back-propagation and gradient descent could be applied to improve compositional hierarchies, while learning of convolutional networks could be augmented with the explicit spatial structure from compositional hierarchies to allow to learn from principles of compositional hierarchies or to add more control when combining with the structure from other domains, such as using depth information.

The remaining of this paper is structured as follows: in Section II Gaussian parametrization is presented in detail, in Section III experimental evaluation is given and concluding remarks are drawn in Section IV.

## II. GAUSSIAN PARAMETRIZATION

In this section we present our model in detail. We also show direct relation to deep convolutional networks on one side and compositional hierarchical models on the other side.

We first provide notation for deep convolutional neural networks and then derive our filter parametrization from it. In deep neural network each neuron's activation is modeled with a linear function wrapped around a specific non-linearity:

$$y_i = f\left(\sum_s w_s \cdot x_s + b_s\right),$$

with output activation  $y$ , input activations  $x$ , bias  $b$  and weights  $w$  determining linear combination of input activations further modified by a non-linear function  $f(\cdot)$ . Since we are limiting to image domain neuron's activations are organized as 3-dimensional matrix  $N \times M \times S$ , where two dimensions,  $N \times M$ , represent image or feature plane, and the third dimension,  $S$ , represent different channels. With the introduction of weight-sharing along 2D feature plane the convolutional neural network has the model of neuron's activation map  $Y$  defined as:

$$Y_i = f\left(\sum_s W_s \odot X_s + b_s\right),$$

where  $\odot$  is convolution operation applied between  $X_s$  as  $N \times M$  activation map from  $s$ -th channel and  $W_s$  as  $A \times B$  filter for  $s$ -th channel. This gives  $N - A/2 \times M - B/2$  output map for each channel and element-wise summation over channels represents final activation map  $Y_i$  after non-linearity  $f(\cdot)$  is applied as well. Typically, several activation maps  $Y_i$  are created, and network is organized in layers such that  $X_s$  in  $l$ -th layer is output activation  $Y_i$  from  $l - 1$  layer.

We propose to further parametrize filter weights  $W_s$  using a mixture of Gaussian distributions:

$$W_s = \sum_k \tilde{w}_k G(\vec{\mu}_k, \sigma_k), \quad (1)$$

where filter weights are composed from  $K$  number of Gaussian components  $G(\vec{\mu}_k, \sigma_k)$ , each with  $\vec{\mu}_k$  mean,  $\sigma_k^2$  variance and  $\tilde{w}_k$  weight factor. Parameters for Gaussian components are different for each filter channel  $s$ , but we omit this in notation. It stems from two dimensional filter  $W_s$  of size  $A \times B$  that  $G(\vec{\mu}, \sigma)$  is two dimensional matrix of the same size as well:

$$G(\vec{\mu}, \sigma)[a, b] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(a-\mu[0])^2 + (b-\mu[1])^2}{2\sigma^2}}.$$

We use two dimensional means but single dimensional variance for simplification. Commonly used filter sizes are  $3 \times 3$ ,  $5 \times 5$  or  $9 \times 9$ ; however, such small sizes can lead to significant discretization errors in  $G(\vec{\mu}_k, \sigma_k)$ . We avoid this by replacing normalization factor computed in continuous space with one computed in discretized space leading to our final distribution function  $G(\vec{\mu}, \sigma)$ :

$$G(\vec{\mu}, \sigma)[a, b] = \frac{1}{N(\vec{\mu}, \sigma)} g(a, b, \vec{\mu}, \sigma),$$

where  $g(a, b, \vec{\mu}, \sigma)$  is non-normalized Gaussian distribution and  $N(\vec{\mu}, \sigma)$  is a sum over non-normalized Gaussian distribution computed for a filter of size  $A \times B$ :

$$g(a, b, \vec{\mu}, \sigma) = e^{-\frac{(a-\mu[0])^2 + (b-\mu[1])^2}{2\sigma^2}}, \\ N(\vec{\mu}, \sigma) = \sum_a \sum_b g(a, b, \vec{\mu}, \sigma).$$

Note, that proposed parametrization is similar to Gaussian mixture model; however, in our case we do not enforce  $\sum \tilde{w} = 1$  and component weights can take any value,  $\tilde{w} \in [-\infty, \infty]$ . Having negative weights is important to approximate edges as differences between neighboring components, where positive components can be interpreted as requirement for a presence of a feature and negative components as requirements for an absence of a feature. Nevertheless, this could be also satisfied by normalizing to a sum of absolute components weights,  $\sum |\tilde{w}| = 1$ , but this further complicates gradient computation and in our case did not help much.

Learning of filter values as weights  $W_s$  is now replaced with the learning of several Gaussian components each defined with parameters weight  $\tilde{w}_k$ , mean  $\vec{\mu}_k$  and variance  $\sigma_k^2$ . The number of Gaussian components per filter can be considered as hyper-parameter. Learning can be performed the same as in

convolutional networks with the gradient descend and back-propagation. Parameters are optimized by computing gradients w.r.t. the cost function  $C(y, \bar{y})$ , which leads to three different types of gradients. By applying chain rule we can define gradient for component weight  $\partial C / \partial \tilde{w}_k$  as a dot-product of back-propagated error and input feature  $X_s$  convolved with  $k$ -th Gaussian component:

$$\frac{\partial C}{\partial \tilde{w}_k} = \sum_{n,m} \frac{\partial C}{\partial Z} \cdot \frac{\partial Z}{\partial \tilde{w}_k} = \sum_{n,m} \frac{\partial C}{\partial Z} \cdot \sum_{a,b} X_s \odot G(\tilde{\mu}_k, \sigma_k), \quad (2)$$

where  $Z = \sum W_s \odot X_s + b_s$  and  $\partial C / \partial Z$  is back-propagated error. Note, that only  $s$ -th channel of input features are used since weight component  $\tilde{w}_k$  appears only in  $W_s$ . Back-propagated error is computed the same as in non-parametrized convolutional network:

$$\frac{\partial C}{\partial Z_s^l} = \frac{\partial C}{\partial Z_s^{l+1}} \odot \text{rot}(W_s),$$

where back-propagated error from previous layer is convolved with the  $180^\circ$  rotated weight filter  $\text{rot}(W_s)$  and weight filter itself can be computed from Eq. 1. We can similarly apply chain rule to obtain gradient for mean and variance:

$$\frac{\partial C}{\partial \mu_k} = \sum_{n,m} \frac{\partial C}{\partial Z} \cdot \sum_{a,b} X_s \odot \frac{\partial G(\tilde{\mu}_k, \sigma_k)}{\partial \mu_k}, \quad (3)$$

$$\frac{\partial C}{\partial \sigma_k} = \sum_{n,m} \frac{\partial C}{\partial Z} \cdot \sum_{a,b} X_s \odot \frac{\partial G(\tilde{\mu}_k, \sigma_k)}{\partial \sigma_k}, \quad (4)$$

where filter derivatives of Gaussian distribution w.r.t. the mean and variance are:

$$\frac{\partial G(\tilde{\mu}_k, \sigma_k)}{\partial \mu_k} = \tilde{w}_k \frac{N(\tilde{\mu}_k, \sigma_k) \cdot \frac{g(\tilde{\mu}_k, \sigma_k)}{\partial \mu_k} - g(\tilde{\mu}_k, \sigma_k) \cdot \frac{\partial N(\tilde{\mu}_k, \sigma_k)}{\partial \mu_k}}{[N(\tilde{\mu}_k, \sigma_k)]^2},$$

$$\frac{\partial G(\tilde{\mu}_k, \sigma_k)}{\partial \sigma_k} = \tilde{w}_k \frac{N(\tilde{\mu}_k, \sigma_k) \cdot \frac{g(\tilde{\mu}_k, \sigma_k)}{\partial \sigma_k} - g(\tilde{\mu}_k, \sigma_k) \cdot \frac{\partial N(\tilde{\mu}_k, \sigma_k)}{\partial \sigma_k}}{[N(\tilde{\mu}_k, \sigma_k)]^2}.$$

Derivatives for mean and variance are also plotted in Figure 1, with first two images depicting derivatives for mean  $\tilde{\mu}_k$  in two spatial dimensions, and the last image depicting derivative for standard deviation  $\sigma_k$ . Filters can be interpreted as edge detectors for mean positions and as blob detector for standard deviation. In the context of gradient learning filters point to the interpretation of Gaussian components moving away from edges and preferring a more blob-like features.

#### A. Implementation

Upper definition for computing forward and backward pass can naturally be implemented by utilizing existing computational models for neural networks. Forward pass in particular can be performed with existing forward pass implementations in regular convolutional networks, with the only exception, that filters  $W_s$  are generated from Gaussian components. This has to be performed at every change of any of the components, therefore at every iteration. The computational requirements

for generating filter depends on several factors: number of components, number of sub-features, number of features and the size of filter. First three factors may lead to explosion and increase computational needs; however, by limiting the size of filter this issue can be significantly reduced. Backward pass can also be performed with existing implementation, since the difference is in added convolution of input data with different Gaussian filters. Existing implementations must already perform convolution operation therefore adding new convolutions would not present any significant implementation challenges. Note, that from Eq. 3 and 4 gradients for mean and variance appear complicated; however, this is in fact nothing more than a filter, since  $G(\tilde{\mu}_k, \sigma_k)$  itself is a filter. As such they can be precomputed for each iteration and computation reduces to implementing the same algorithm as for Eq. 2.

There are also other possible implementation approaches. Instead of looking at the Eq. 1 as a single filter we can decompose it to every single component as individual filters. Normally, this would be inefficient since we know that there would be several components per each filter, and it would be better to run convolution only once. However, when there are many features per layer each one will have its own filter for given sub-feature  $X_s$ . In regular convolutional network convolution with each filter has to be performed separately, since there is no common structure in filters to exploit. But in our model we know that the structure will always be a mixture of Gaussian. Furthermore, the only difference between different components will be mean position and variance. Different position of filters would easily be possible to collapse to a single filter, since different positions can be obtained by just shifting the output from a convolution with a non-offsetted component. While variance can be discretized to several common values, e.g. step 0.1 between ranges of  $[0.5, 5]$  would yield 45 filters, making this implementation efficient when more than 50 features per layers are used. Note, that there are possible limitations to this implementation. For instance, offsetting the output instead of component's position would only work for discretized positions, while discretizing variance would not be possible with multidimensional variance. In this work we do not focus on evaluating different computational implementations, but only point to their possibility, and use the simple extension of existing convolutional networks in our evaluations.

#### B. Deep convolutional networks as specialization

Deep convolutional networks can now be considered as a special case of parametrized convolutional network. We can initialize Gaussian model with  $N \times M$  Dirac deltas. If Dirac deltas are evenly spread over the filter such that each one exclusively covers only one pixel, then a single component weight  $\tilde{w}_k$  in Eq. 1 would correspond to a single value in filter weights  $W_s$ . In our model we can approximate each Dirac delta with a Gaussian component with a small, near-zero variance. Gradients of mean and variance would, in this case, collapse to zero since derivative of Dirac delta is zero, while gradient of weight would simplify to:

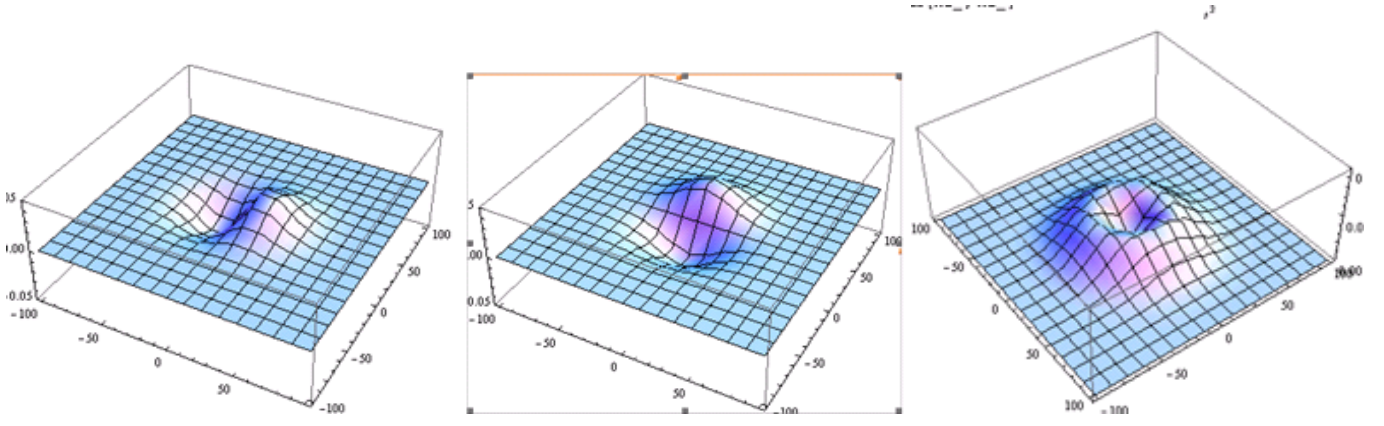


Figure 1. Derivatives for mean  $\bar{\mu}_k$  in two dimension and standard deviation  $\sigma_k$  used in Gaussian component's gradient for gradient descent learning.

$$\frac{\partial C}{\partial \tilde{w}_k} = \sum_{n,m} \frac{\partial C}{\partial Z} \cdot X_s[\bar{\mu}_k],$$

which is an equation for weight gradient in standard convolutional neural network.

### C. Neural network as compositions

Gaussian parametrization of weight opens the door for compositional interpretation of the network. An output activation map  $Y_i$  can be considered as a response map to a presence or an absence of a particular part at different image locations. That part is defined as a composition of several sub-parts who's presence is recursively defined by the previous layer. A response to a particular sub-part corresponds to input activation map  $X_s$ , and how a set of sub-parts are composed into a part is governed by corresponding weight filters  $W_s$ . Since convolution and correlation are related operations a weight filter  $W_s$  defines the blueprint for position of sub-part relative to the center of the filter. Furthermore, since weight filter  $W_s$  is composed from Gaussian component, they can now be considered as a model governing a particular sub-part. Mean  $\bar{\mu}_k$  defines the offset from the center of the part where the sub-part is searched for, variance  $\sigma_k^2$  defines how much a sub-part is allowed to vary in 2D spatial dimensions, and weight  $\tilde{w}_k$  defines the importance of the sub-part. In particular, higher positive weight  $\tilde{w}_k$  value correlates strongly with the presence of the sub-part, while higher negative value correlates strongly with the absence of the sub-part. Small, near-zero values do not correlate with neither presence nor absence, and indicate that sub-part is not important for the feature.

The above interpretation is missing a key element of neural networks, non-linearity. Depending on non-linearity employed different interpretations can be used. With the sigmoid function a response map is mapped into values between  $[0, 1]$  and can be interpreted as a probability map of a specific feature. While ReLU can be considered as a thresholding of parts. If part has response above zero, then they are passed forward and considered as being found with a specific response value, while if they are below zero, they are discarded and considered not present at all. Bias  $b_s$  can additionally modify this behavior,

and move the threshold from zero to any value deemed appropriate for training data.

Compositional interpretation makes the most sense when weights are sparse. Each part feature  $Y_i$  is always composed from *all* sub-part features from the previous layer. Having high weights  $\tilde{w}_k$  for all sub-parts would correspond to a part with many compositions. No meaningful semantical interpretation may be given when there are 100+ sub-features. However, a key element in making deep neural networks work properly is the regularization of weights. This is implemented as weight decay and enforces a sparse solution of weights. Such solution would have a better semantical interpretation as only a handful of sub-parts would be required by each part.

## III. EVALUATION

In this section we evaluate the viability of proposed method. We analyze it on a small synthetic toy case where we can control what kind of filters the network needs to learn. Furthermore, we apply it on a classification problem by evaluating it on CIFAR-10 and PaCMan datasets.

### A. Toy/synthetic case

Gaussian parametrization introduced two new types of parameters that have to be optimized, mean and variance. In Eq. 3 and 4 we showed that gradients can be computed for both of them, and gradient descent can theoretically be applied to optimize them. However, equations show that gradients depend on weight, mean and variance, indicating that a proper setting for hyper-parameters, in particular learning rate, will be important to make the gradient descent work in practice. In this section we analyze optimization of our model with the gradient descent method and show which hyper-parameters are important to achieve proper convergence.

All experiments in this section are performed on synthetic models. In particular, we generate a set of filters representing the groundtruth model that our method would need to learn. We assume that the groundtruth model is composed from several Gaussian components. In classification this may not necessary hold true, but we are currently analyzing how hyper-parameters for gradient descent need to be designed to properly converge, and want to avoid adding additional

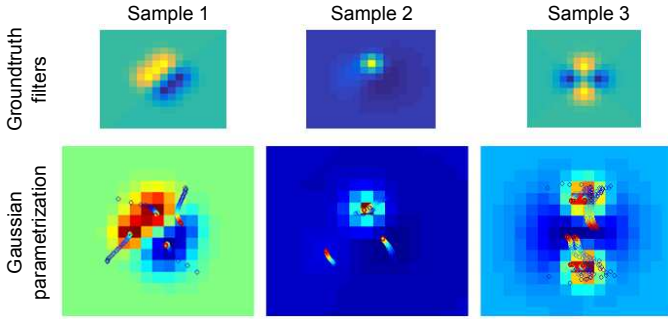


Figure 2.

noise. We use generated filters to compose groundtruth convolutional network. A forward pass on a sample image with the groundtruth network would represent the desired output, and cross-entropy loss between desired output and output of the model being optimized represents the cost function. Note, there is not a single output activation neuron, but multiple neurons each corresponding to specific input image's pixel location. This holds true only because we always use stride of one and no max-pooling between the layers. All filters are of size  $17 \times 17$  giving enough space for components to move around without falling outside of the filter. We use 1000 images from CIFAR-10 dataset as training samples and track the difference between groundtruth filters and trained filters as the measurement of convergence.

*Single layer:* We first experimented with a single layer and a single feature, meaning, input data has only one channel, i.e., gray-scale input image, and only one first-layer feature. Three examples of groundtruth filters are shown in first row in Figure X. The second row shows optimized Gaussian components after 10 epochs. In all three samples small blue/red circles indicate the movement of components, with initial positions in blue, and final positions in red. Clearly, optimizing position with gradient descent converges to correct solution in this simple case. Variances in all three cases seems to be accurately learned as well; however, there are noticeable deviations in negative components in the third sample.

Adding additional sub-features, i.e., using 3-channel image instead of gray-scale, shows similar results which can be seen in Figure X. Note, that in our case it was important to avoid RGB input image and use Lab color-space instead to get correct results. This proved to be important for regular CNNs as well, since RGB channels are highly correlated, and cause ambiguity when back-propagating the error.

It is also worth noting that gradient descent is sensitive to incorrect initializations. To achieve the above results we initialized component weights with normal distribution and positioned components by evenly spreading them over the filter. We also started with small standard deviations of 0.5.

*Two layers:* Next, we experimented with two-layer network. We used two first-layer features composed from three sub-features, corresponding to 3-channel input data. On the second layer we used one feature composed from two sub-features. The results of optimization are shown in Figure X. There are several important details that allow to achieve such results. It proved important to slightly over-parametrize the problem

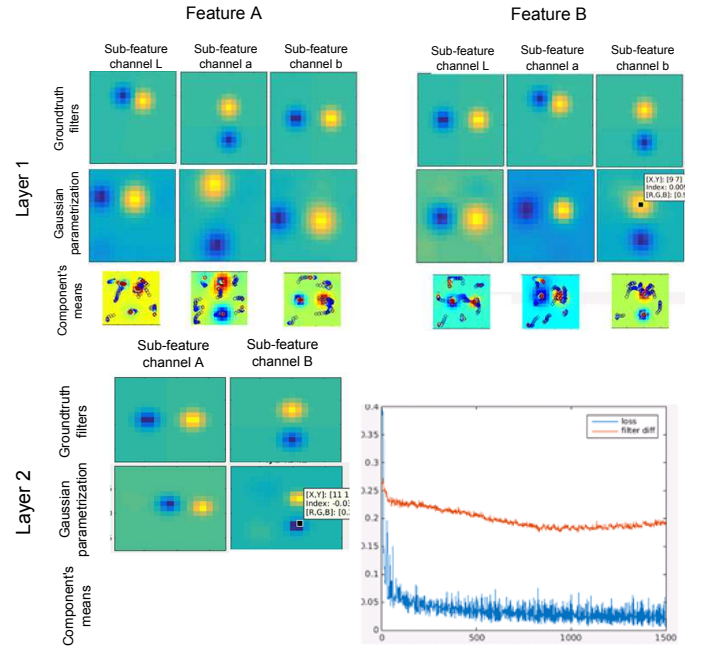


Figure 3.

by using  $3 \times 3$  or  $5 \times 5$  Gaussian components. Components also need to be evenly spread over the filter. If this is ensured gradient descent will quickly push component weights to at least have correct weight signs at appropriate positions. Wrong weights can quickly have negative effect on mean and variance gradients. Despite of this, plain gradient descent may still not be able to achieve desired result. It most cases the variance was incorrect and many times collapsed close to zero. Position of some components was correct in some cases, but stalled in other cases due to collapse of variance to zero. We were able to partially mitigate this problem by adjusting learning rates for means and variance; however, it turned out learning rate should be adjusted on per-component basis, since adjusting learning rate helped in some examples, but did the opposite in others. We achieved proper per-component adjustment by utilizing AdaDelta with a small amount of momentum and weight decay. AdaDelta adjusts learning rate based on history of individual gradient, thus allowing to smooth-out the gradient over time and prevent them from overshooting if higher learning rate needs to be used for other gradients.

The experiments also showed that when there are more Gaussian components available than are needed, multiple components will be used to describe the same area of the filter. This effect can be seen in Figure Y. This may not be the desired behavior since components are not efficiently used. We experimented with adding the penalty function that forces Gaussian components not to overlap. This can be implemented as intersection area between two components which penalizes the solutions with significant overlaps. However, this proved not to work as expected, since now penalty function prevented components from converging to correct position. We therefore did not opt to address this issue this way, but it may be possible to address it later by, e.g., merging components close to each other or randomly moving only one of them

	num features	stride	CNN with Gaussian parametrization		non-parametrized CNN
			filter size	num components	filter size
conv1	32	1	$7 \times 7$	$2 \times 2$	$3 \times 3$
relu1	/	1	/	/	/
pool1	/	1	$3 \times 3$	/	/
conv2	32	1	$9 \times 9$	$3 \times 3$	$5 \times 5$
relu2	/	1	/	/	/
pool2	/	2	$3 \times 3$	/	/
ip1	10	/	$15 \times 15$	/	$15 \times 15$

Table I

	num features	stride	CNN with Gaussian parametrization		non-parametrized CNN
			filter size	num components	filter size
conv1	16	1	$9 \times 9$	$2 \times 2$	$5 \times 5$
relu1	/	1	/	/	/
pool1	/	2	$3 \times 3$	/	/
conv2	16	1	$9 \times 9$	$2 \times 2$	$5 \times 5$
relu2	/	1	/	/	/
pool2	/	2	$3 \times 3$	/	/
ip1	23	/	$15 \times 15$	/	$15 \times 15$

Table II

away. Nevertheless, this issue does not affect the performance characteristics, since the difference is only in solution with the smaller description.

Note, that regular CNNs achieve results that appear much worse. Figure X shows the model being learned by regular CNN. It is immediately visible that filters being learned are only a crude approximation of groundtruth filters. Nevertheless, the optimization correctly converges on a training set which should be an indicator of severe overfitting. This should be mitigated with a stronger levels of weight decay to force a more sparse solution.

### B. Classification on CIFAR-10

Next, we applied the model to the classification problem by evaluating it on CIFAR-10 dataset. Evaluation is performed with a network containing three layers. First two layers were convolutional and the third one is fully-connected. We use soft-max with multinomial logistic loss as cost function. Detailed network configuration is shown in Table X. Input data was three-channel RGB data additionally mean-subtracted, but non-normalized to unit variance, since between each layer we use ReLU non-itineraries which are less sensitive to input variance. Note, we have slightly bigger filters size in parametrized CNN, but use less components. We also restrict component's positions and standard deviation to ensure derivative of filters will not have many non-zero values fall outside of the valid filter size and also preventing them from stalling if collapsing to a single point. Positions were restricted to at least 1.5 pixels away from filter borders, and standard deviations were restricted to above 0.5.

Figure X shows results when training CNN with Gaussian parametrization after 50 epochs. Comparing with non-parametrized CNN shown in Figure Y, we can achieve the same accuracy with approximately the same speed of learning. Both achieve slightly less than 70% accuracy on testing set. An example of a few filters are shown in Figure.X...

### C. Classification on PaCMan-DB

We evaluated similar network on PaCMan dataset as well. Dataset contains gray-scale and depth images generated from 3D models of 23 categories of various kitchen objects, with each category containing between 15 and 30 samples. Each object is captured at dense, regular viewpoint intervals, but we use only 28 different viewpoints, summing to a total of around 40.000 images. Note, we use only gray-scale images. Dataset

was split into around 20.000 samples for training and 20.000 for testing. We ensured all viewpoint images of one sample are in the same split and each category has proportionally the same number of objects in testing and training. Input images are resized to  $128 \times 96$ . Network configuration was slightly modified to accommodate higher resolution images. Detailed network configuration is shown in Table Y.

...

## IV. CONCLUSION

....

## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *Advances In Neural Information Processing Systems* 25, 2012, pp. 1097–1105.
- [2] M. D. Zeiler and R. Fergus, "Visualizing and Understanding Convolutional Networks," in *European Conference on Computer Vision*, nov 2013, pp. 818–833. [Online]. Available: <http://arxiv.org/abs/1311.2901>
- [3] A. Mahendran and A. Vedaldi, "Understanding Deep Image Representations by Inverting Them," in *Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [4] Y. LeCun, "Generalization and network design strategies," *Connectionism in perspective*, pp. 143–155, 1989.