
Gaussian process classifiers and CNN uncertainty

Sebastian Borgeaud dit Avocat
spb61@cam.ac.uk

LE49 - Probabilistic Machine Learning Project

Abstract

The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Introduction

Model uncertainty is often of crucial importance.

Softmax output does not give model confidence, maybe show the example used by Yarin Gal.

With model uncertainty it is possible to treat uncertain inputs and special cases explicitly. Example if model uncertainty is high we might decide to pass the input to a human for classification.

Another example is in Reinforcement learning: with uncertainty a RL agent can decide when to explore and when to exploit in an environment. With principled uncertainty for the agent's Q-value function, it can learn much faster using techniques such as Thompson sampling.

Yarin Gal shows how dropout in neural networks can be interpreted as a Bayesian approximation of a Gaussian Process.

1.1 Convolutional neural networks [1] (ADD MORE CITATIONS)

Convolutional neural networks (CNNs) were originally inspired by biological processes. They have become standard in many deep learning applications, especially in image processing or vision tasks. A convolutional neural network is a type of feedforward neural network, typically consisting of convolution layers, pooling layers and fully connected layers:

- **Convolution layers** are composed of several convolution kernels each computing a different feature map. The output feature maps are obtained by convolving the input with the convolution kernel and then applying an element-wise nonlinearity. Mathematically, the feature value $z_{i,j,k}^l$ at location (i, j) of the k^{th} feature map in the l^{th} layer is computed as:

$$z_{i,j,k}^l = \mathbf{w}_k^l \mathbf{x}_{i,j}^l + b_k^l$$

where \mathbf{w}_k^l and b_k^l are the weight and bias vectors for the k^{th} convolution kernel in the l^{th} layer and $\mathbf{x}_{i,j}^l$ is the input patch centered around (i, j) in the l^{th} layer. The output value is computed by applying a nonlinearity $a(\cdot)$ point-wise:

$$x_{i,j,k}^{(l+1)} = a(z_{i,j,k}^l)$$

- **Pooling layers** aim to achieve shift-invariance and reduce the number of parameters in the network by reducing the resolution of the feature maps. The pooling layer operates on each

feature map independently. Mathematically, the output of a pooling layer with pooling operation $\text{pool}(\cdot)$ is given by

$$y_{i,j,k}^l = \text{pool}(x_{m,n,k}^l), \forall (m, n) \in \mathcal{R}_{i,j}$$

where $\mathcal{R}_{i,j}$ is a local neighbourhood around (i, j) . Typically, the pooling operation computes the average or the maximum.

- **Fully connected layers** connect every neuron in the previous layer to single neuron in the current layer. Fully connected layers are the layers used in standard neural networks. Mathematically, the output of a fully connected layer is given by:

$$x_i^{(l+1)} = a\left(\left(\sum_j w_{i,j}^l x_j^l\right) + b_i^l\right)$$

where $a(\cdot)$ is a nonlinearity, $w_{i,j}^l$ is the weight connecting neuron j in the l^{th} layer to neuron i in layer $(l+1)^{\text{th}}$, and b_i^l is the bias weight for neuron i .

Optimising a convolutional neural network is done in the same way as optimising standard neural networks. A differentiable loss function is computed for the training examples (often done in batches) and the gradients w.r.t. the weights of the network are computed. Using these gradients, the weights are updated in a gradient descent step. Typically, more complex update rule that take into account change momentum (e.g. Adam [2]) are used as they converge faster to a local minimum.

1.2 Gaussian processes

Formally, a Gaussian Process is defined as a collection of random variables, any finite number of which have (consistent) joint Gaussian distributions. A Gaussian process therefore defines a distribution over functions and is fully specified by a mean function $m(x)$ and a covariance function $k(x, x')$. Write $f \sim \mathcal{GP}(m, k)$ meaning f is distributed as a GP with mean m and covariance k .

Using the GP we can draw samples from the function for any finite number n of locations. Given locations $\mathbf{x} = [x_1, \dots, x_n]$, first compute $\mu_i = m(x_i)$, $\Sigma_{i,j} = k(x_i, x_j)$. We can then sample a vector from this distribution: $\mathbf{f} \sim \mathcal{N}(\mu, \Sigma)$.

1.2.1 Regression

We can now use this GP as a prior for Bayesian inference. Let \mathbf{f} be the known function values for the training examples and let \mathbf{f}_* be the set of function values corresponding to the set of test input X_* . The joint distribution is given by

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} = \mathcal{N}\left(\begin{bmatrix} \mu \\ \mu_* \end{bmatrix}, \begin{bmatrix} \Sigma & \Sigma_* \\ \Sigma_*^T & \Sigma_{**} \end{bmatrix}\right)$$

where μ_* are the test means, Σ_* are the training-test covariances, and Σ_{**} are the test-test covariances. Since we know the training values \mathbf{f} , we are interested in the conditional distribution of \mathbf{f}_* given \mathbf{f} :

$$\mathbf{f} | \mathbf{f}_* \sim \mathcal{N}(\mu_* + \Sigma_*^T \Sigma^{-1}(\mathbf{f} - \mu), \Sigma_{**} - \Sigma_*^T \Sigma^{-1} \Sigma_*)$$

This corresponds to a posterior Gaussian process, given by

$$\begin{aligned} f | \mathcal{D} &\sim \mathcal{GP}(m_{\mathcal{D}}, k_{\mathcal{D}}), \\ m_{\mathcal{D}}(x) &= m(x) + \Sigma(X, x)^T \Sigma^{-1}(\mathbf{f} - \mathbf{m}) \\ k_{\mathcal{D}}(x, x') &= k(x, x') - \Sigma(X, x)^T \Sigma^{-1} \Sigma(X, x') \end{aligned}$$

where $\Sigma(X, x)$ is a vector of covariances between every training case in X and x . Furthermore, it is easy to incorporate noise in the observations. Assuming i.i.d. additive Gaussian noise, every $f(x)$ now has extra covariance with itself with a magnitude equal to the noise variance σ_n^2 :

$$f | \mathcal{D} \sim \mathcal{GP}(m_{\mathcal{D}}, k_{\mathcal{D}} + \delta_{ii} \sigma_n^2)$$

where $\delta_{ii'} = 1$ iff $i = i'$ is the Kronecker's delta.

The mean function $m(x)$ and the covariance function $k(x, x')$ are typically parametrised in terms of hyper-parameters θ . During training we find the values of the hyper-parameters which optimise the marginal likelihood:

$$L = \log p(\mathbf{y}|\mathbf{x}, \theta) = -\frac{1}{2}|\Sigma| - \frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{y} - \boldsymbol{\mu}) - \frac{n}{2} \log(2\pi)$$

This optimisation can be done using standard gradient methods.

1.2.2 Classification

Binary classification using Gaussian processes can be done by setting a GP prior over a latent function $f(\mathbf{x})$ and then using squashing function such as the sigmoid to obtain a probability:

$$\pi(\mathbf{x}) = p(y = +1|\mathbf{x}) = \sigma(f(\mathbf{x}))$$

Inference is done in two steps. First, compute the distribution of the latent variable corresponding to a new test input \mathbf{x}_*

$$p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int p(f_*|\mathbf{X}, \mathbf{x}_*, \mathbf{f})p(\mathbf{f}|\mathbf{X}, \mathbf{y})d\mathbf{f}.$$

Second, use this distribution to compute a probabilistic prediction

$$\bar{\pi}_* = p(y_* = +1|\mathbf{X}, \mathbf{y}, \mathbf{x}_*) = \int \sigma(f_*)p(f_*|\mathbf{X}, \mathbf{y}, \mathbf{x}_*)df_*$$

As the likelihood is non longer Gaussian, the first integral becomes analytically intractable. Similarly, depending on the sigmoid function, the second integral can also be intractable. Hence, we need to use approximations, either analytical or numerical, for example using Monte Carlo sampling, to solve the integrals.

Multi-class classification is typically [3] approached by assuming the following labelling rule for y_* given \mathbf{x}_* :

$$y_* = \arg \max_k f^k(\mathbf{x}_*), \text{ for } k = 1, \dots, C$$

where each $f^k(\cdot)$ is a nonlinear latent function with a GP prior. The likelihood is again non-Gaussian meaning that approximation techniques have to be used to perform inference and to optimise the hyper-parameters.

TODO: Explain how classification is done with GPs, explain the difficulties. Say something about approximation techniques?

1.3 Uncertainty in Deep Learning

Obtaining uncertainty bounds in deep learning would be without doubt of great benefit to many fields. An autonomous car could make much safer decisions knowing how much uncertainty there is in the consequence of taking different actions. When using current deep learning methods we only have point estimates of parameters and predictions (CITE DL BOOK). Despite the fact the output of a softmax layer in a classification task can be interpreted as a probability distribution $p(y_i = k) = f(x_i)_k$, it doesn't say anything about how certain the model is: The model could assign a high probability to a class but can still be highly uncertain about it (CITE GAL).

One way to get uncertainty from deep learning models was discovered by Yarin Gal (GAL THESIS). In fact, there is a rather subtle but deep connection between Gaussian Process and neural networks. In particular, he showed that using Dropout not only at training time but also at test time, and obtaining multiple estimates for a single test example allows us to compute the uncertainty in a principled manner by taking the sample variance of those predictions (CITATION).

In this project, I focus on a different way to obtain those uncertainty bounds: Using a Gaussian process trained on the features extracted by a CNN that was trained in a previous step. It should be noted that this is a less theoretically grounded approach than the one taking by Gal (CITE), but nonetheless provides a practical way of getting uncertainty whilst keeping the power of neural networks.

2 Method

2.1 Convolutional Neural Network architecture

The first step consists in training a convolutional neural network on the MNIST dataset, which can be done using one of the many deep learning libraries. For example, I use Keras (CITATION) which defines an extra abstraction layer above TensorFlow. The network architecture is provided in the Keras tutorial for image classification on MNIST. The first two layers are convolutional layers 3×3 kernels and ReLU activations, where $\text{ReLU}(x) = \max(x, 0)$. The first layer has 32 feature maps and the second layer has 64. A max-pooling layer with kernel size 2×2 is then applied to the output of the convolutional layer. The final 2 layers are fully connected layers, with hidden sizes of 128 and 10 respectively. The first fully connected layer has a ReLU activation. The last fully connected layer uses a softmax activation, which outputs a probability distribution over the 10 classes representing the 10 digits. Furthermore, Dropout (CITATION) is applied after the max-pooling layer with $p = 0.25$ and after the first fully connected layer with $p = 0.5$, where p is the probability of dropping a neuron. The network is trained using an Adadelta optimiser over 10 epochs with batches of size 128.

2.2 Gaussian process

Using the trained CNN model, we can extract the activations in the last hidden layer. These correspond to a vector of size 128 for each input image, and can be thought of as the features extracted by the CNN for classification. I then train a Gaussian Process using these features. This has the advantage of being a much lower dimensional classification task compared to classifying the entire input image directly, which has dimension $28 \times 28 = 784$. However, the MNIST dataset, with 60,000 training images, is still considered a large dataset for GPs because inference takes $\mathcal{O}(n^3)$ time where n is the number of training instances [4]. To make the training possible, I use GPFlow [5], which implements various approximation algorithms for Gaussian Processes. Furthermore, the library is built on top of TensorFlow which has the further advantage of being usable on a GPU out-of-the-box, which provides a further speed-up. More precisely, I use the Sparse Variational Gaussian Process Classifier presented by Hensman et al. [4]. The model learns a set of inducing points of size m which are then used instead of the training points. Typically, $m < n$, which makes the task tractable as its complexity is $\mathcal{O}(nm^2)$. For the experiments done in this project I chose to use 600 inducing points, which allowed the model to be trained in about 15 minutes on a NVidia Titan Xp GPU.

I train the Gaussian process with a simple kernel as it is sufficient to obtain good accuracy. The kernel consists of a Matern kernel with $\nu = \frac{3}{2}$ (CITATION) and a white noise function to account for the noise in the input data.

TODO: Add figure of the model...

2.3 Evaluation

2.3.1 Accuracy

The standard metric for evaluating a model f on a classification task is the accuracy, which can be computed as follows:

$$\text{accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(f(\mathbf{x}^{(i)}) = y^{(i)})$$

where N is the number of test data points and $(\mathbf{x}^{(i)}, y^{(i)})$ is the i^{th} test point. Note that the accuracy does not take into account the uncertainty with which a model makes the prediction and only rewards correctly classified examples.

2.3.2 Reject option

One way to incorporate the model uncertainty into an evaluation metric is by allowing the model to say “I don’t know”. For example, if the model was presented with an image of the letter ‘a’, it would still classify the image into one of the 10 digits, which might not be a desired outcome. Instead, the model should be able to reject classification for the inputs in which the uncertainty is too high. This is known as **classification with a reject option** [6]. An appropriate metric assigns a cost α to

misclassification and a cost β to the reject option \textcircled{R} :

$$\text{R-accuracy}_{\alpha,\beta} = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(f(\mathbf{x}^{(i)}) = y^{(i)}) - \alpha \mathbb{1}(f(\mathbf{x}^{(i)}) \neq y^{(i)}) - \beta \mathbb{1}(f(\mathbf{x}^{(i)}) = \textcircled{R})$$

Here, α and β are parameters that are set by the user. A large value for α corresponds to cases where misclassification is expensive, for example in medical diagnosis or when a self-driving car has to make a steering decision. The β parameter intuitively indicates how much value the user places on obtaining a prediction, as a higher β will increase the cost of the reject option.

Now, given uncertainty bounds with each class probability, we can design the model to reject classification according to certain rejection rules. For example, we could decide to reject if the standard deviation of the highest probability is above a pre-defined threshold. Another rejection rule could decide to reject if the second highest class probability lies within a certain factor ϵ of the standard deviations of the highest probability, i.e. the model would reject if

$$|f_p(\mathbf{x}^{(i)})_1 - f_p(\mathbf{x}^{(i)})_2| < \epsilon f_{\text{std}}(\mathbf{x}^{(i)})_1$$

where $f_p(\mathbf{x}^{(i)})$ is the classification probability vector for input $\mathbf{x}^{(i)}$ with $f_p(\mathbf{x}^{(i)})_j$ being defined as the probability of the j^{th} most likely class, and $f_{\text{std}}(\mathbf{x}^{(i)})_1$ is the standard deviation of the most likely class returned by the model for input $\mathbf{x}^{(i)}$. Note that the ϵ parameter defines a trade-off between the number of misclassified examples at the number of examples that would have been correctly classified but were rejected instead. For the results presented next, I focus on this decision rule with different values for ϵ .

3 Results

3.1 MNIST dataset

The MNIST dataset is considered a toy dataset in deep learning and is generally seen as a “unit test” for models. In fact, using the model provided in the Keras tutorial and presented above, I obtain an accuracy of 99.12% on the held-out test images. This means that only 88 of the 10,000 test images are misclassified.

Using the Gaussian process with the CNN features, the accuracy of the model on the test images improves to 99.28%, reducing the number of misclassified images to 72. Although this increase in performance might seem surprising at first, it can be explained by the fact that the Gaussian process is more complex and more powerful than the last layer of the CNN, which just computes a linear combination of the features and then applies the softmax activation.

| Model | Accuracy | Misclassified images |
|-------|----------|----------------------|
| CNN | 99.12% | 88 |
| GP | 99.28% | 72 |

Table 1: MNIST test images accuracy for the CNN and the GP.

Furthermore, it is interesting to inspect the uncertainty the Gaussian process assigns to different test images. In order to do this, I plot the cumulative distribution of standard deviations for the most likely class for correctly and incorrectly classified test images in figure 1. For the correctly classified examples, the model’s uncertainty is low: the model gives the correct class probability with a standard deviation of less than 0.1 for 97.3% of the correctly classified images. However, for the incorrectly classified images, the standard deviation of the most likely class seems to be higher: The model predicts a standard deviation below 0.1 for only 19.4% of the incorrectly classified images. This suggests that the model is generally less certain of its prediction in the cases where it got the wrong class.

3.2 Noisy MNIST (n-MNIST)

Next, I evaluate both the CNN and the Gaussian process on the n-MNIST dataset [7]. The dataset consists of the MNIST images, but modified in 3 different ways so as to make the classification task

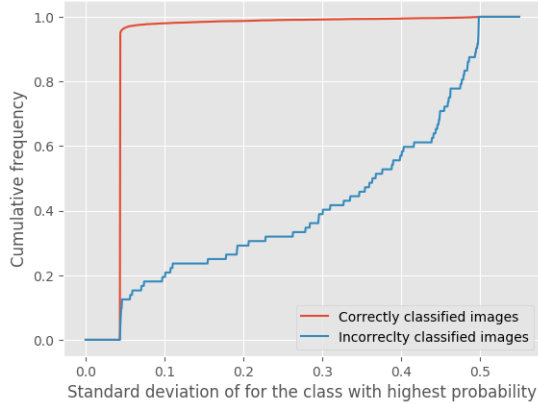


Figure 1: Cumulative distribution of standard deviations for correctly and incorrectly classified test images. The standard deviation of the most likely class is taken for both correctly and incorrectly classified images.

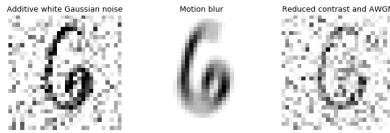


Figure 2: Example digits from the noisy-MNIST dataset. From left to right: Additive white Gaussian noise, Motion blur, and Reduced contrast with AWGN

harder. The first set of images is created by adding white Gaussian noise with a signal-to-noise ratio of 9.5 to the images. The second set of images is obtained by emulating a linear motion blur. The last set of images is the hardest to classify as it is created by reducing the contrast and also adding white Gaussian noise to the original images. Figure 2 shows an example image taken from the 3 datasets.

The task of classifying examples whose underlying distribution is different from the underlying distribution of the training images is called **out-of-distribution** classification (CITE GAL DISSERTATION). The performance of the two models is shown in table 2. Although the GP performs better on the NMIST test images, the CNN performs better on all 3 noisy datasets. This could again be explained by the fact that the GP is a more complex model, therefore is overfitting more on the MNIST data and hence is not able to generalise as well on these different datasets. Especially noticeable is the difference in performance for the reduced contrast dataset, where the CNN obtains an accuracy of 77.71% compared to 70.13% for the GP. As the 3rd dataset is the least similar the original one, this higher discrepancy supports the claim that the GP is overfitting more on the MNIST dataset.

| Model | Accuracy | | |
|-------|----------|-------------|------------------|
| | AWGN | Motion blur | Reduced contrast |
| CNN | 94.35% | 93.27% | 77.71% |
| GP | 93.34% | 92.09% | 70.13% |

Table 2: MNIST test images accuracy for the CNN and the GP.

We can again inspect how uncertain the model is when classifying the noisy dataset by plotting the cumulative distributions of the standard deviations of the most likely digit in of the 3 datasets. Figure 3 shows the resulting plots, which indeed suggest that the uncertainty increases as the model is asked to classify images that are increasingly different from those it was trained on. The uncertainty for the additive white Gaussian noise dataset is higher than the uncertainty for the standard MNIST test set. This is shown by the fact that the cumulative distribution for the correctly classified images in the AWGN images is lower. Similarly, the cumulative distribution further decreases for the motion

blur and reduced contrast images, with a lower classification accuracy corresponding to a higher uncertainty.

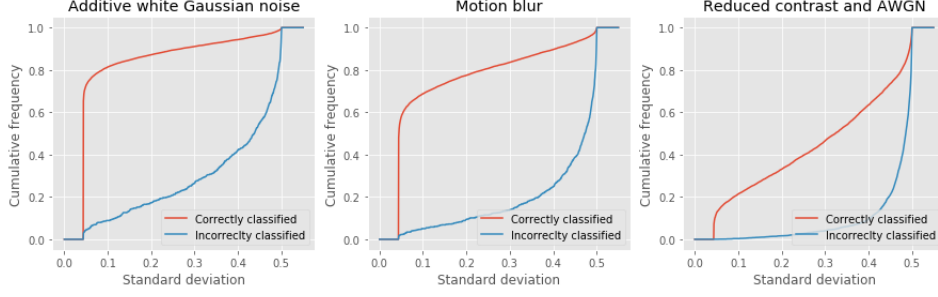


Figure 3: Cumulative distribution of standard deviations for correctly and incorrectly classified test images. The standard deviation of the most likely class is taken for both correctly and incorrectly classified images.

3.3 Classifying with a reject option

Finally, I also evaluate the performance of the GP augmented with a reject option according to the rejection rule described in 2.3.2. Figure 4 (left) shows the R-accuracy with varying α and fixed $\beta = 0$ for different values of ϵ . In the cases where $\alpha > 1$, that is when the cost of misclassification is higher than the reward obtained by a correct classification, there seem to be an optimal value of ϵ lying between 20 and 40. On the right, the number of misclassified but accepted and the number of correctly classified but rejected examples is plotted for varying ϵ . The plot shows that there is indeed a trade-off between these two quantities and that the trade-off decision can be made by choosing an appropriate ϵ .

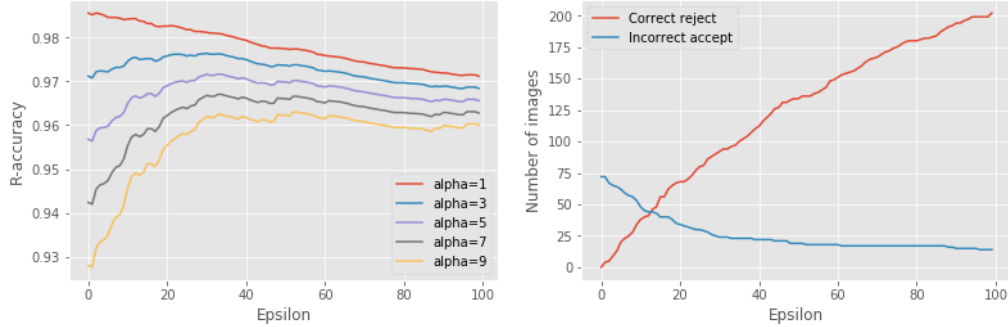


Figure 4: Right: R-accuracy obtained using varying ϵ and α but setting $\beta = 0$. For $\alpha > 0$ there seems to be an optimal value of ϵ around 30. Left: Number of examples misclassified or rejected. This illustrates the trade-off between the number of misclassified but accepted and correctly classified but rejected examples.

4 Discussion

Only 72 of the 10,000 MNIST test images are misclassified by the GP. To gain some insights into why only those were misclassified, the 72 misclassified input images are shown in figure 8 in appendix A. The images reveal that many of those 72 images are indeed very ambiguous, even for a human annotator. For example, it is hard to say whether the first image should be a ‘1’ or a ‘2’. This also explains why the uncertainty given by the model is higher for the misclassified images than it is for the correctly classified images. In fact, one could argue that for those ambiguous images the model should indeed be less certain of its prediction.

TODO: Discuss definition of the metrics.

4.1 Uncertainty with additive white noise

A further property the model uncertainty should obey is that as an image becomes more and more distorted, the uncertainty should increase. To verify that this holds I add increasingly more white Gaussian noise to a test digit and then ask the GP to make predictions for those images. The results are shown in figure 5 for an image of the digit ‘6’. Note that for the last few images, it becomes

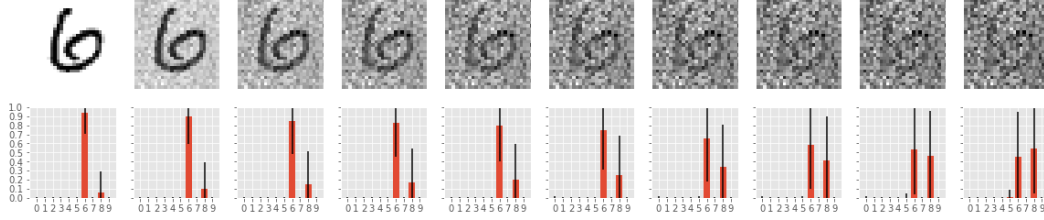


Figure 5: Top: Test images with increasingly more additive white Gaussian noise. Bottom: Gaussian process predictions for the test images. The error bars show the standard error for each class.

almost impossible to discern the original digit from the noise (especially when looked at individually, rather than next to the original image). Indeed, the model wrongly classifies the last image as a ‘5’, although does so with a high uncertainty. Figure 6 shows the probability and uncertainty assigned to the digit ‘6’ for the 10 images. As hoped for, the uncertainty increases as more noise is added to the image. Especially noteworthy is the increase of the standard deviation relative to the class probability.

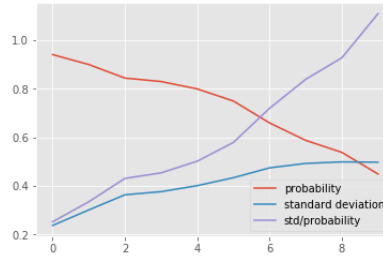


Figure 6: Probability and standard deviation assigned to class ‘6’ for the 10 images shown in figure 5. Also plotted is the standard deviation relative to the class probability.

4.2 CNN features visualisation

Finally, I show the features learned by the CNN by reducing the number of dimensions of the features to 2 using the t-SNE [8] dimensionality reduction technique. The learned features are shown for the MNIST test dataset and for the additive white Gaussian noise images of the n-MNIST dataset.

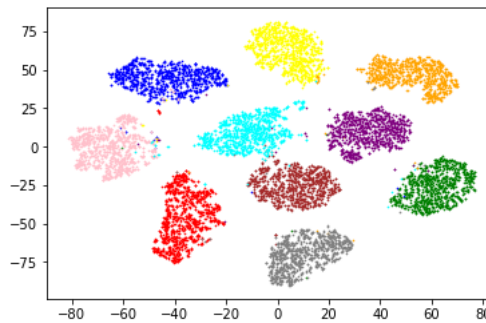


Figure 7: Test point embedding using t-SNE

References

- [1] Jiuxiang Gu, Zhenhua Wang, Jason Kuen, Lianyang Ma, Amir Shahroudy, Bing Shuai, Ting Liu, Xingxing Wang, and Gang Wang. Recent advances in convolutional neural networks. *CoRR*, abs/1512.07108, 2015.
- [2] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [3] Carlos Villacampa-Calvo and Daniel Hernández-Lobato. Scalable multi-class gaussian process classification using expectation propagation. *arXiv preprint arXiv:1706.07258*, 2017.
- [4] James Hensman, Alexander G de G Matthews, and Zoubin Ghahramani. Scalable variational gaussian process classification. 2015.
- [5] Alexander G. de G. Matthews, Mark van der Wilk, Tom Nickson, Keisuke Fujii, Alexis Boukouvalas, Pablo Le'on-Villagr'a, Zoubin Ghahramani, and James Hensman.
- [6] C Chow. On optimum recognition error and reject tradeoff. *IEEE Transactions on information theory*, 16(1):41–46, 1970.
- [7] Saikat Basu, Manohar Karki, Sangram Ganguly, Robert DiBiano, Supratik Mukhopadhyay, and Ramakrishna R. Nemani. Learning sparse feature representations using probabilistic quadrees and deep belief nets. *CoRR*, abs/1509.03413, 2015.
- [8] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

A Misclassified MNIST images by the Gaussian process model

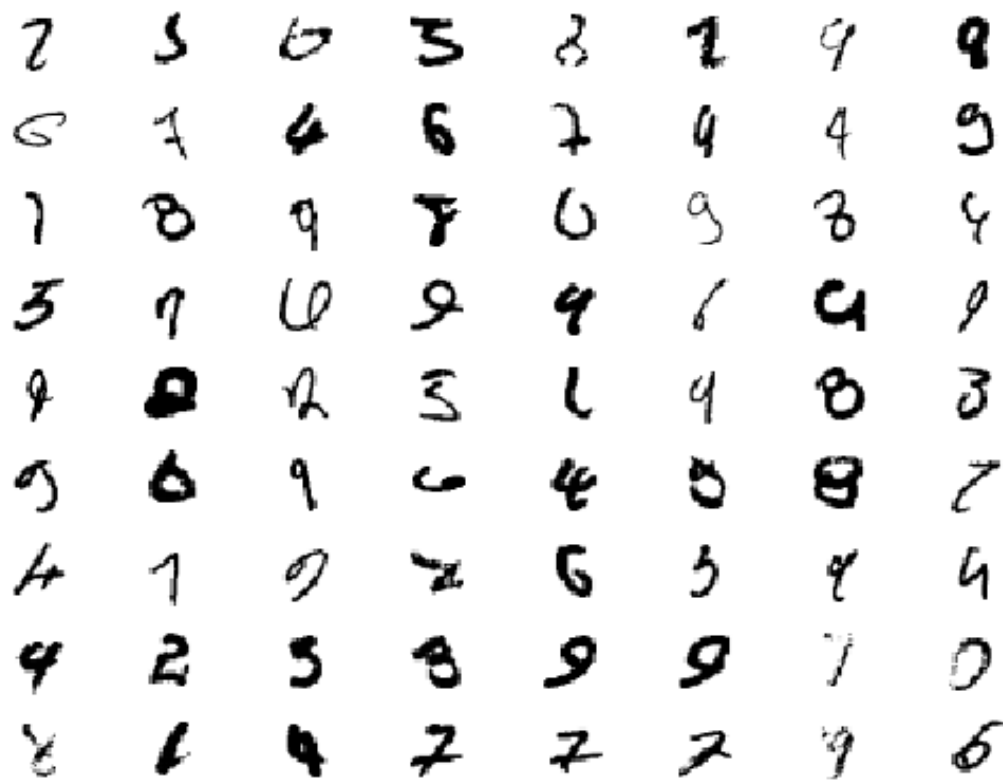


Figure 8: Test images misclassified images by the GP model.