
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

2 Method

2.1 Convolutional Neural Network architecture

The first step consists in training a convolutional neural network on the MNIST dataset. This is made very easy through different libraries, for example Keras (CITATION) which can be used as an extra abstraction layer above TensorFlow. I use the network architecture provided in the Keras tutorial for image classification on MNIST. The network consists of two convolutional layers with 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 layer uses a softmax, 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 Adadelat optimiser over 10 epochs with batches of size 128.

2.2 Gaussian process

2.3 Metrics

2.3.1 Accuracy

2.3.2 0-d-1 loss?

2.3.3 With reject cost

3 Results

4 Discussion

More quantitative evaluation of the results. Show examples of images that are missclassified. Discuss definition of the metrics.

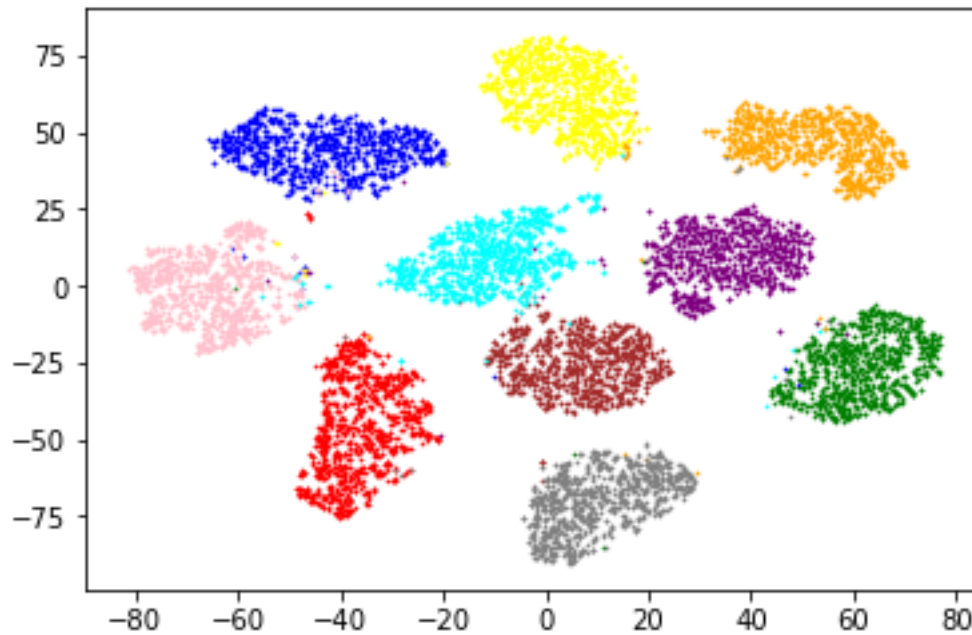


Figure 1: Test point embedding using t-SNE

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.