

Bayesian deep learning

Deep Learning Practical Work

Aymeric DELEFOSSE & Charles VIN

2023 – 2024



Contents

1	Bayesian Linear Regression	2
1.1	Linear Basis function model	2
1.1.1	Gaussian basis functions	2
1.2	Non Linear models	5
1.2.1	Polynomial basis functions	5
1.2.2	Gaussian basis functions	6
2	Approximate Inference in Classification	9
2.1	Bayesian Logistic Regression	9
2.1.1	Maximum-A-Posteriori Estimate	9
2.1.2	Laplace Approximation	9
2.1.3	Variational Inference	11
2.2	Bayesian Neural Networks	12
2.2.1	Variational Inference with Bayesian Neural Networks	12
2.2.2	Monte Carlo Dropout	13
3	Uncertainty Applications	15
3.1	Monte-Carlo Dropout on MNIST	15
3.2	Failure prediction	16
3.3	Out-of-distribution detection	17

Chapter 1

Bayesian Linear Regression

1.1 Linear Basis function model

1.1.1 Gaussian basis functions

1.2. Recall closed form of the posterior distribution in linear case. Then, code and visualize posterior sampling. What can you observe? Let N denote the number of training examples, K the number of dimensions of the outputs, and p the number of features (or predictors) in the input data. Consider $X \in \mathbb{R}^{N \times p}$, the input matrix, and $Y \in \mathbb{R}^{N \times K}$, the output matrix.

We typically define the posterior distribution as $p(w|X, Y)$. This distribution represents our updated beliefs about the parameters w after observing the data X and Y . According to Bayes' rule, we can express the posterior distribution as the product of two components:

1. $p(w)$, which represents our prior beliefs about the distribution of w before observing the data.
2. $p(Y|X, w)$, indicating the probability of observing the data Y given the parameters w and the data X . This quantifies how likely our data is under different hypotheses represented by w .

Therefore, the formula for the posterior distribution is given by:

$$p(w|X, Y) \propto p(Y|X, w)p(w)$$

In this formula, we start with our initial beliefs (priors) and then update these beliefs based on the new data (likelihood). In the case of our linear model, we know that $y_i = \Phi_i^T w + \epsilon$, with $\Phi \in \mathbb{R}^{N \times (p+1)}$ representing the design matrix and ϵ denoting the residual. Assuming that the error follows a centered Gaussian distribution with standard deviation $\beta^{-1} = 2\sigma^2$, meaning that $\epsilon \sim \mathcal{N}(0, \beta^{-1})$. Consequently, we can conclude that:

$$p(y_i|x_i, w) \sim \mathcal{N}(\Phi_i^T w, \beta^{-1})$$

Furthermore, we selected a centered Gaussian prior with a variance of $\alpha^{-1}I$ where α governs the prior distribution over the weights w :

$$p(w|\alpha) \sim \mathcal{N}(0, \alpha^{-1}I)$$

In this specific case, we can demonstrate that the posterior distribution $p(w|X, Y)$ follows a Gaussian distribution as follows:

$$p(w|X, Y) \sim \mathcal{N}(\mu, \Sigma)$$

The precision matrix Σ , which is the inverse of the covariance matrix of the distribution parameters, is defined as:

$$\Sigma^{-1} = \alpha I + \beta \Phi^T \Phi$$

The mean of the distribution parameters is given by:

$$\mu = \beta \Sigma \Phi^T Y$$

Parameters α and β serve analogous roles, with α governing the prior distribution and β regulating the likelihood.

Now, we can proceed to sample from the updated (with data) posterior distribution $p(w|X, Y)$. This process is demonstrated in Figure 1.1, with $\alpha = 2$ and $\beta = (2 \times 0.2^2)^{-1}$. It is worth noting that when $N = 0$, the posterior distribution $p(w|X, Y)$ simplifies to the prior distribution $p(w)$. As we increase the number of

data points, we can observe how the model's certainty increase, demonstrated by the reduction in the variance of the parameters of the distribution. In other words, having more data points reduces aleatoric uncertainty.

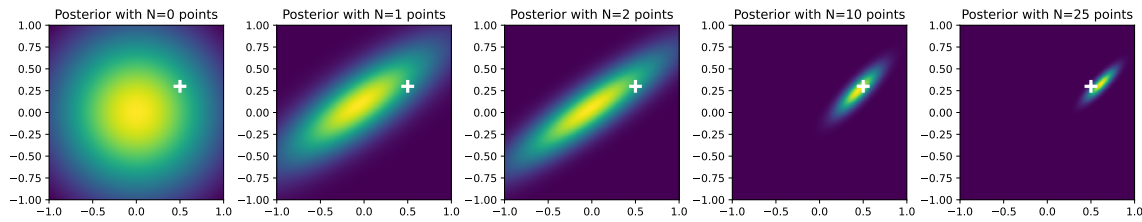


Figure 1.1: Evolution of a Bayesian posterior distribution with increasing data points: A visual representation of Bayesian inference, depicting how the posterior distribution updates as more data is incorporated. From left to right, the figures show the posterior with $N = 0$ (prior distribution), $N = 1$, $N = 2$, $N = 10$, and $N = 25$ data points, respectively. The cross mark represents the ground truth.

1.3. Recall and code closed form of the predictive distribution in linear case. Using the information from the previous question, we can now calculate the predictive distribution for new data point x^* by marginalizing over the parameter w . This is represented by the following equation where $\mathcal{D} = \{X, Y\}$ denotes the dataset:

$$p(y|x^*, \mathcal{D}, \alpha, \beta) = \int p(y|x^*, w, \beta) p(w|\mathcal{D}, \alpha, \beta)$$

By leveraging the property that the convolution of two Gaussian distributions results in another Gaussian distribution, we can demonstrate that the closed-form expression for the predictive distribution in the linear case is as follows:

$$p(y|x^*; \mathcal{D}, \alpha, \beta) = \mathcal{N}\left(y; \mu^T \Phi(x^*), \frac{1}{\beta} + \Phi(x^*)^T \Sigma \Phi(x^*)\right)$$

We can see that the variance in the predictive distribution σ_{pred}^2 for a new observation can be divided into two parts:

1. The aleatoric uncertainty, which represents the inherent noise in the data, that we fixed around β^{-1} ;
2. The epistemic uncertainty related to the model parameters w , characterized by $\Phi(x^*)^T \Sigma \Phi(x^*)$.

It's worth noting that as the number of data points N approaches infinity ($\lim_{N \rightarrow \infty} \Phi(x^*)^T \Sigma \Phi(x^*) = 0$), our understanding of the model parameters becomes nearly perfect. In this scenario, the only remaining source of uncertainty is the aleatoric uncertainty, stemming from the noise in the data.

1.4. Based on previously defined `f_pred()`, predict on the test dataset. Then visualize results using `plot_results()` defined at the beginning of the notebook. Figure 1.2 illustrates a comparison between the predictions made by a Bayesian Linear Regression model and the actual ground truth. In the left panel, you can see the model's linear fit to the training data, shown as blue points, alongside the true ground truth represented by the green line. To visualize the model's predictive uncertainty, shaded areas are used, ranging from dark to light, which correspond to one, two, and three standard deviation intervals, respectively.

The right panel of the figure focuses on the predictive variance σ_{pred}^2 along the x-axis. This variance is depicted by a curve that widens as it moves away from the center of the training data, marked by the vertical dashed line. These visual elements together provide a comprehensive view of the model's confidence in its predictions across the entire domain.

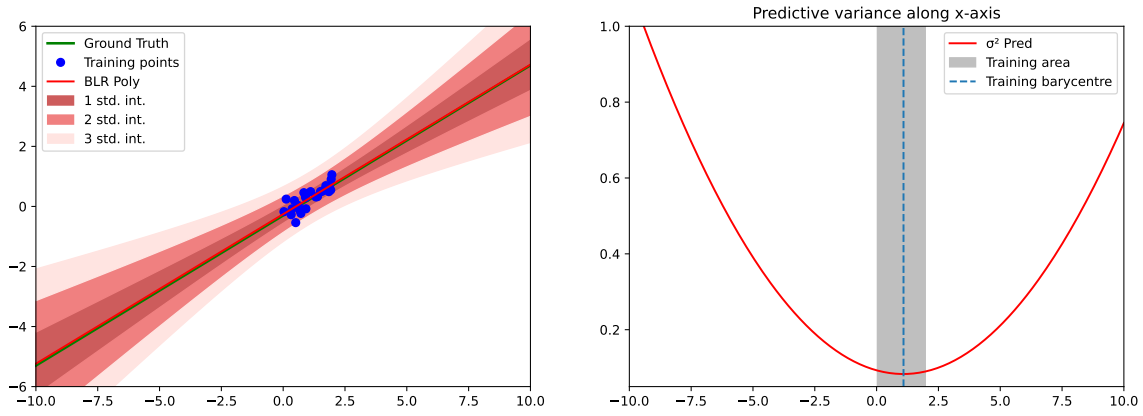


Figure 1.2: Visualization of predictive distribution of a linear dataset using Bayesian Linear Regression. The left panel illustrates the fit (red line) to the training data (blue points) against the ground truth (green line). The shaded areas represent the predictive uncertainty, with one, two, and three standard deviation intervals shown in progressively lighter shades. The right panel displays the predictive variance σ_{pred}^2 across the x-axis. The vertical dashed line indicates the center of the training data.

1.5. Analyse these results. Why predictive variance increases far from training distribution? Prove it analytically in the case where $\alpha = 0$ and $\beta = 1$. In the left panel of Figure 1.2, we can observe that the confidence intervals are narrowest near the cluster of training points. This suggests higher confidence in predictions within this area. As we move away from the center of the training data (towards the extremities of the x-axis), the confidence intervals become wider, indicating increasing uncertainty in the model's predictions.

Looking at the right panel of Figure 1.2, we notice that the variance remains low in the region where the training data is located (the grey shaded area). This correlates with the tight confidence intervals shown in the left panel. As expected, the predictive variance is lowest near the training barycentre, reflecting greater model certainty (i.e. lower epistemic uncertainty) in this region due to the presence of more training data points. As we move away from the training area on either side, the predictive variance increases significantly, which is consistent with the expanding confidence intervals in the left panel. This sharp increase in variance indicates a significant decrease in the model's confidence (i.e. a significant increase in the model's epistemic uncertainty) in its predictions outside the range of the training data. This pattern is expected, as the model has less points to rely on for predictions in these regions.

Let's prove it analytically. With $\alpha = 0$ and $\beta = 1$, the computation of Σ^{-1} simplifies as follows:

$$\begin{aligned}\Sigma^{-1} &= 0 \cdot I_3 + 1 \cdot \Phi^T \Phi \\ &= \Phi^T \Phi \\ &= \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_N \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \\ &= \begin{pmatrix} N & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}\end{aligned}$$

To invert Σ , we use the classic formula for a 2×2 matrix:

$$\Sigma = \frac{1}{\det \Sigma^{-1}} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & N \end{pmatrix}$$

Returning to our expression for $\sigma_{\text{pred}}^2 = \Phi(x^*)^T \Sigma \Phi(x^*)$, we have:

$$\begin{aligned}
 \sigma_{\text{pred}}^2 &= \Phi(x^*)^T \Sigma \Phi(x^*) = \frac{1}{\det \Sigma^{-1}} \begin{pmatrix} 1 \\ x^* \end{pmatrix} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & N \end{pmatrix} \begin{pmatrix} 1 \\ x^* \end{pmatrix} \\
 &= \frac{1}{\det \Sigma^{-1}} \begin{pmatrix} \sum x_i^2 - x^* \sum x_i \\ x^* N - \sum x_i \end{pmatrix} \begin{pmatrix} 1 \\ x^* \end{pmatrix} \\
 &= \frac{1}{\det \Sigma^{-1}} \left(\sum_{i=1}^N x_i^2 - x^* \sum_{i=1}^N x_i + x^{*2} N - x^* \sum_{i=1}^N x_i \right) \\
 &= \frac{1}{\det \Sigma^{-1}} \left(\sum_{i=1}^N x_i^2 - 2x^* \sum_{i=1}^N x_i + x^{*2} N \right) \\
 &= \frac{1}{\det \Sigma^{-1}} \sum_{i=1}^N (x_i - x^*)^2
 \end{aligned}$$

This formula reveals that the predictive variance is directly proportional to epistemic uncertainty, which, in the case of our linear regression, is the squared differences between each training data point x_i and the prediction point x^* . As x^* moves away from the region where the training data is concentrated, these squared differences grow, consequently leading to an increase in the predictive variance.

Bonus: What happens when applying Bayesian Linear Regression on the following dataset? Examining the right panel in Figure 1.3, an intriguing observation is made: the variance is unexpectedly minimized at the barycenter, i.e. the "hole" where there are no training data points. Normally, one would anticipate an increase in variance in data-scarce regions. However, this phenomenon can be explained by the fact that a point within this region is positioned closely between two clusters of training points. This proximity results in lower squared differences (epistemic uncertainty) and, as previously discussed, leads to a lower predictive variance. Additionally, due to the broader distribution of our training dataset, we observe lower predictive variance at the endpoints when compared to our previous results.

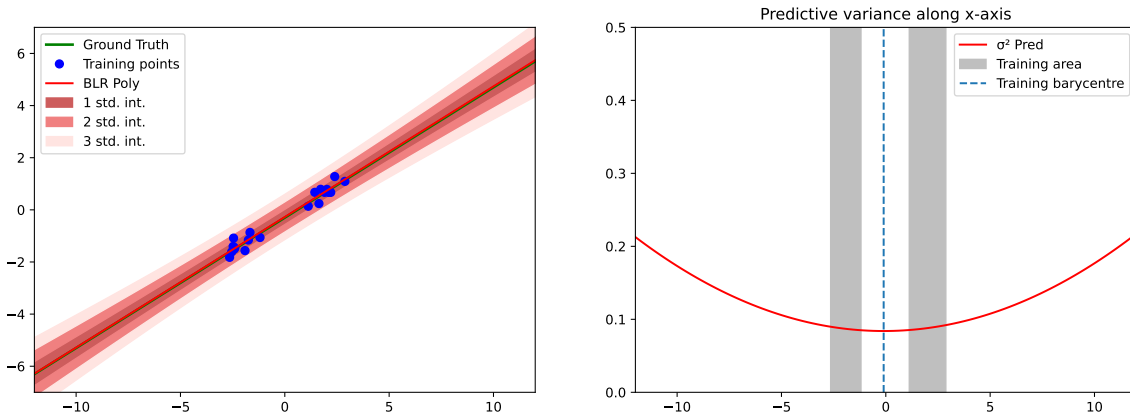


Figure 1.3: Visualization of predictive distribution of a linear dataset featuring a "hole" using Bayesian Linear Regression.

1.2 Non Linear models

1.2.1 Polynomial basis functions

2.2. Code and visualize results on sinusoidal dataset using polynomial basis functions. What can you say about the predictive variance? Figures 1.4 and 1.5 illustrate a comparison between the predictions made by a Bayesian Polynomial Regression model and the actual ground truth. Given that the closed-form solution for the posterior and predictive distribution in Φ is similar to the linear case, we can draw the same conclusions: as we move further away from our training data, uncertainty increases. However, due to our polynomial kernel, the model exhibits higher predictive variance (i.e. higher epistemic uncertainty) when considering the values alone. Nevertheless, thanks to this basis function, our model demonstrates the ability to capture more intricate patterns, such as those resembling a sinusoidal function. It's worth noting that beyond the training data points, the model struggles to generalize, as it lacks the necessary data points to accurately capture the periodic nature of the sinusoidal ground truth.

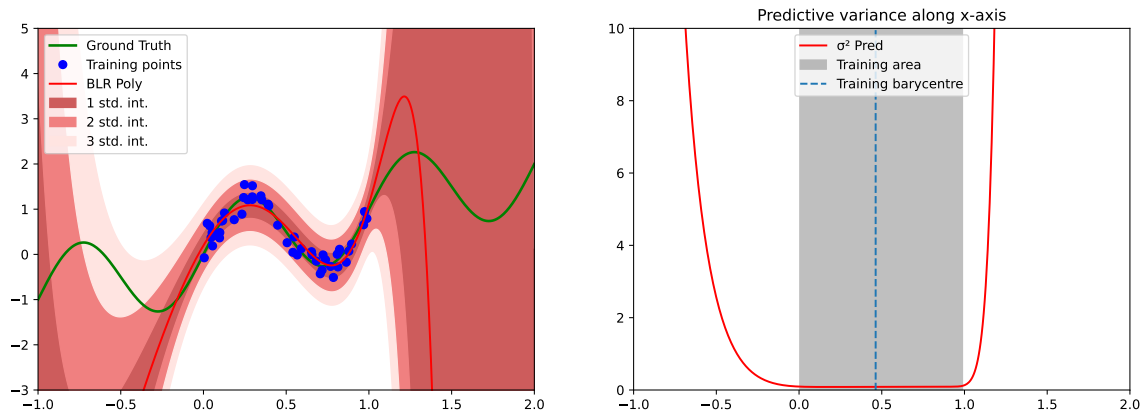


Figure 1.4: Visualization of predictive distribution of a sinusoidal dataset using Bayesian Polynomial Regression.

Furthermore, we decided to explore the effects of providing data points that are more dispersed and less abundant, as shown in Figure 1.5. Unsurprisingly, we observe that the variance around these small clusters is minimal, and as we move away from these regions, the variance increases. This phenomenon highlights the strength of these models: when given data, they become increasingly confident in specific areas. Therefore, in the context of uncertainty, we can gradually narrow down our desired outcomes by continuously adding more data points over time.

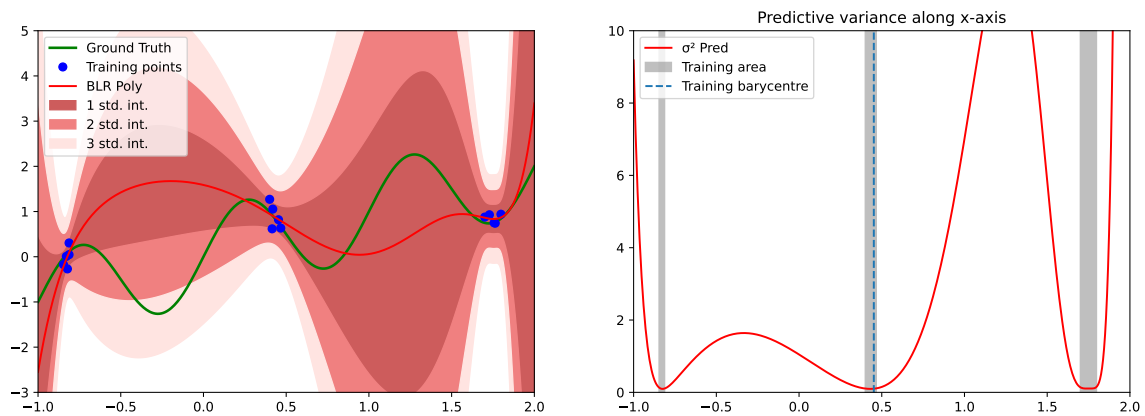


Figure 1.5: Visualization of predictive distribution of a sinusoidal dataset featuring "holes" and sparse data points using Bayesian Polynomial Regression.

1.2.2 Gaussian basis functions

2.4. Code and visualize results on sinusoidal dataset using Gaussian basis functions. What can you say this time about the predictive variance? Figures 1.6 and 1.7 offer a comparison between the predictions generated by an RBF Gaussian Regression model and the actual ground truth. In contrast to the previous two models, our analysis leads to distinct conclusions. Notably, the predictive variance is at its highest not in regions far from the training barycenter. This is because the mean encompasses the entire range of data, resulting in predictive variance being concentrated solely within this zone. Beyond this range, predictive variance diminishes, which we explain in the next question. The predictive variance demonstrates a fluctuating pattern, marked by peaks and troughs corresponding to regions where the density of training points varies. Consequently, the model closely adheres to the mean of the training data, which is evident in its impact on the behavior of the sinusoidal curves.

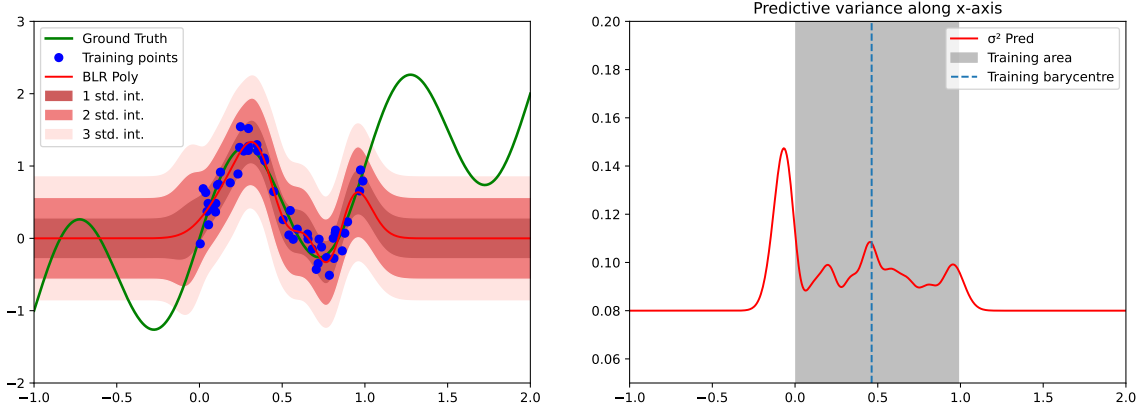


Figure 1.6: Visualization of predictive distribution of a sinusoidal dataset using RBF Gaussian Regression, where we set $\mu \in [0.0, 1.0]$ and $M = 9$.

When we visualize our model with a reduced number of data points, it reinforces what has been discussed and emphasizes the role of hyperparameters. Notably, the model struggles to capture the information provided by the data points at the extremes, as they lie outside the range of the specified mean. Nonetheless, the predictive variance exhibits a significant increase as anticipated in regions where no data points are present.

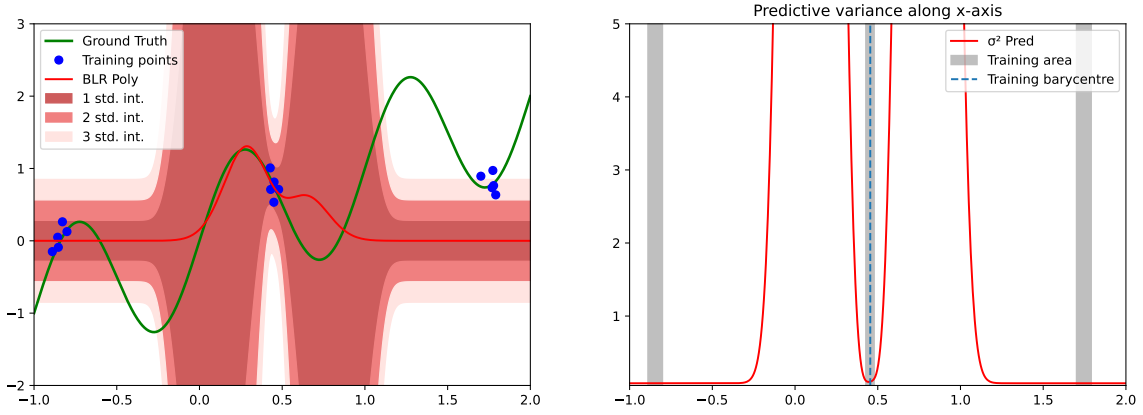


Figure 1.7: Visualization of predictive distribution of a sinusoidal dataset featuring "holes" and sparse data points using RBF Gaussian Regression, where we set $\mu \in [0.0, 1.0]$ and $M = 9$.

2.5. Explain why in regions far from training distribution, the predictive variance converges to this value when using localized basis functions such as Gaussians. Let us recall the definition of the Radial Basis Function:

$$\Phi_j(x_i) = \exp\left(-\frac{(x_i - \mu_j)^2}{2s^2}\right),$$

where μ_j represents the center of the j -th Gaussian basis function and s is a parameter controlling the spread of the Gaussian.

The primary observation is that this Gaussian function reaches its maximum at its center, where $x = \mu_j$ since $(x - \mu_j)^2 = 0$. As the distance $(x - \mu_j)^2$ increases, the exponential term rapidly tends toward zero. So, if x^* is far from μ_j , we can approximate $\Phi_j(x^*) \approx 0$. As a result, when we multiply this vector by the covariance matrix and its transpose, the epistemic uncertainty converges toward zero. Therefore, $\sigma_{\text{pred}}^2 = \beta^{-1} + \Phi(x^*)^T \Sigma \Phi(x^*) \approx \beta^{-1}$, signifying that the predictive variance is effectively reduced to aleatoric uncertainty. This can be observed in Figures 1.6 and 1.7, where $\sigma_{\text{pred}}^2 = \beta^{-1} = 0.08$.

Consequently, μ and M are two critical hyperparameters (with s derived from them), and their selection should be based on our data. For instance, setting $\mu \in [-2, 2]$ without changing M yields significantly improved results, as demonstrated in Figures 1.8 and 1.9. We observe that these outcomes outperform the polynomial approach while maintaining similarity, thereby highlighting the robustness and popularity of the RBF kernel.

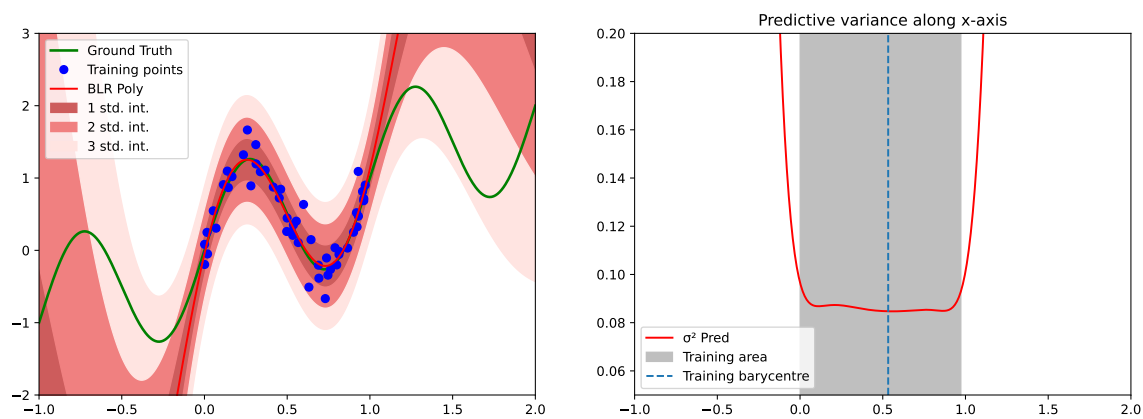


Figure 1.8: Visualization of predictive distribution of a sinusoidal dataset using RBF Gaussian Regression, where we set $\mu \in [-2.0, 2.0]$ and $M = 9$.

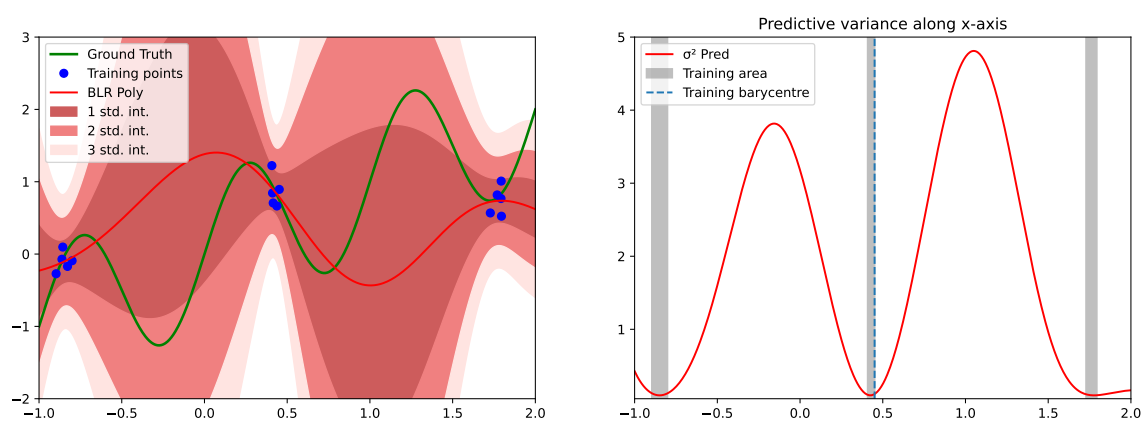


Figure 1.9: Visualization of predictive distribution of a sinusoidal dataset using RBF Gaussian Regression, where we set $\mu \in [-2.0, 2.0]$ and $M = 9$.

Chapter 2

Approximate Inference in Classification

2.1 Bayesian Logistic Regression

2.1.1 Maximum-A-Posteriori Estimate

1.1. Analyze the results provided by Figure 2.1. Looking at $p(y = 1|x, \mathbf{w}_{\text{MAP}})$, what can you say about points far from train distribution? Approximating $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$ with a Dirac delta function is essentially akin to approximating the predictive distribution using \mathbf{w}_{MAP} , meaning $p(y = 1|x, \mathbf{w}_{\text{MAP}}) \approx p(y = 1|x, \mathbf{Y})$. This approximation is quite straightforward.

As depicted in Figure 2.1, the boundary decision remains linear and the model's uncertainty doesn't significantly increase far from the training data. Essentially, it provides a confidence measure for the linear boundary. This indicates that the point-wise estimate of the parameters can only confidently assign points to their respective classes but lacks the capacity to provide nuanced uncertainty measures for points that deviate far from the training data distribution. Thus, this approach is not effective in assessing the uncertainty of outliers or points not well-represented in the training set.

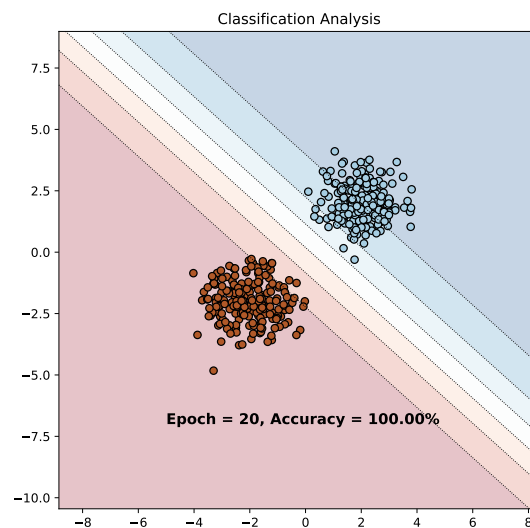


Figure 2.1: Illustration of a Bayesian Logistic Regression model applied to a binary classification task with uncertainty display, with a weight decay of 5×10^{-2} . Two distinct data point clusters — blue and red — represent separate classes, while the surrounding shaded areas reflect the model's predictive uncertainty, with lighter shades indicating lower confidence.

2.1.2 Laplace Approximation

1.2. Analyze the results provided by Figure 2.2. Compared to previous MAP estimate, how does the predictive distribution behave? Compared to the MAP estimate, Bayesian Logistic Regression using the Laplace approximation better captures uncertainty about the model parameters. While the MAP estimate gives a single point estimate of the weights (\mathbf{w}_{MAP}) and therefore a single decision boundary, the Laplace approximation treats the weights as a normal distribution. This distribution is centered around \mathbf{w}_{MAP} and

has a covariance matrix based on the Hessian of the log posterior. This approach allows for uncertainty in the decision boundary, as clearly shown in Figure 2.2.

For instance, moving in the southwest direction away from the red cluster — directly opposite the blue cluster — the model exhibits high confidence that this region predominantly consists of red points. The same holds true for the blue cluster. This directional certainty reflects the model's aleatoric uncertainty, which is the irreducible uncertainty inherent in the observations due to noise or other stochastic effects. Conversely, if we move sideways, out of the line between the clusters to areas with less or no data, the model becomes less certain. This reflects the epistemic uncertainty of our model, relating to what the model doesn't know about its own parameters.

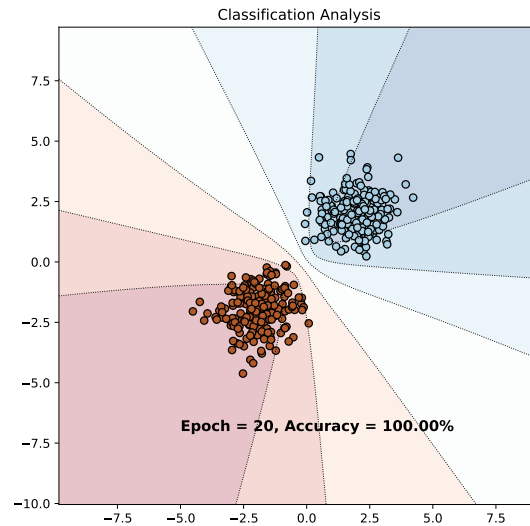


Figure 2.2: Illustration of a Bayesian Logistic Regression with Laplace approximation model applied to a binary classification task, with a weight decay of 5×10^{-2} .

1.3. Comment the effect of the regularisation hyper-parameter WEIGHT_DECAY. The weight decay hyper-parameter controls the complexity of the model. It adds a penalty to the loss function for large weights, effectively encouraging the model to maintain smaller weight values. This usually help prevent overfitting. In Bayesian terms, weight decay corresponds to the precision (inverse variance) of the prior distribution over the weights. A higher weight decay value means a tighter prior, which pulls the weights closer to zero, unless the data provides strong evidence to the contrary. This can affect the predictive distribution by potentially making it more conservative. As a result, the decision boundary may be less flexible and the model may exhibit higher uncertainty, especially in regions far from the training data. This behavior is verified in Figure 2.3. When the weight decay is too high, it results in increased predictive uncertainty (wider shaded areas), whereas when it's too low, it results in high confidence (narrower shaded areas), which may not be justified for unseen data.

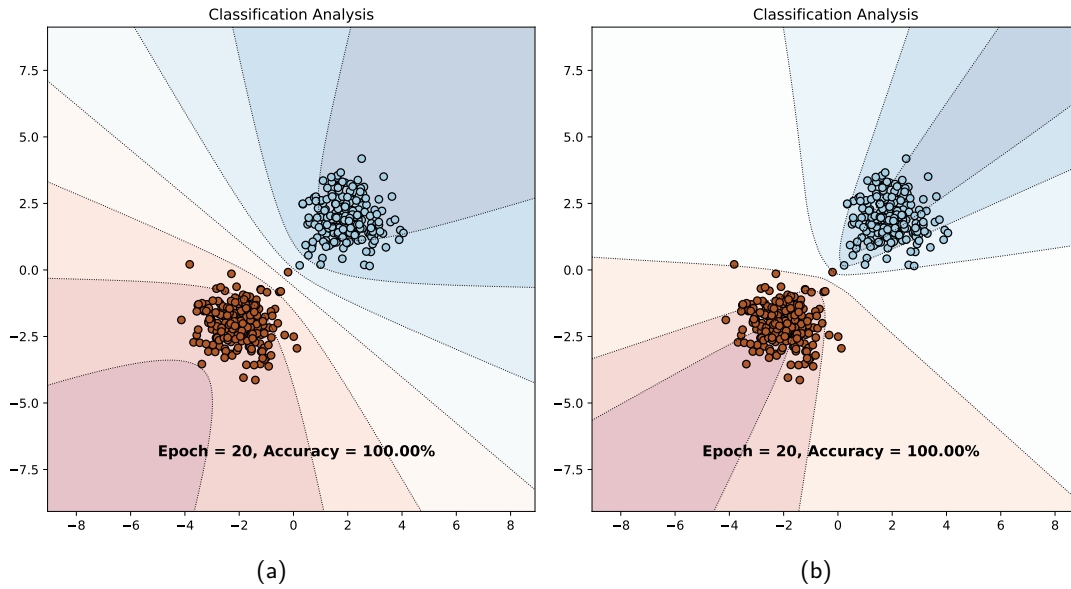


Figure 2.3: Illustration of a Bayesian Logistic Regression with Laplace approximation model applied to a binary classification task, with a weight decay of (a) 0.5 and (b) 5×10^{-5} .

2.1.3 Variational Inference

1.4. Comment the code of the VariationalLogisticRegression and LinearVariational classes. `LinearVariational` represents a single linear layer with variational inference applied. It approximates the weights and biases of the layer with distributions rather than fixed values.

- The class is initialized with the variational parameters for the weights (`w_mu`, `w_rho`) and biases (`b_mu`) of the layer, i.e. the parameters we want to learn. `prior_var` represents the variance of the prior distribution (σ_p^2), which specify our prior belief about the distribution of the weights.
- The `sampling` method uses the reparametrization trick to sample from the variational posterior distribution for the weights $w_i \sim \mathcal{N}(\mu_i, \sigma_i^2)$. To do so, we sample from a centered isotropic multivariate Gaussian where $\sigma^2 = \log(1 + e^p)$ to avoid numerical issues. Thus, $w_i = \mu_i + \sigma_i \odot \epsilon_s$, where $\epsilon_s \sim \mathcal{N}(0, 1)$ is a Gaussian noise. The reparametrization trick allows the gradient of the loss function to backpropagate through the randomness of the sampling process.
- The `kl_divergence` method calculates the Kullback-Leibler divergence between the variational posterior and the prior distribution for the weights:

$$\text{KL}[q_\theta(w) \| p(w)] = \log\left(\frac{\sigma_p}{\sigma_i}\right) + \frac{\sigma_i^2 + \mu_i^2}{2\sigma_p^2} - \frac{1}{2}$$

where σ_p^2 is the variance of our prior distribution $p(w)$ and (μ_i, σ_i^2) the mean and variance of the variational distribution $q_\theta(w)$.

- The `forward` method defines the forward pass by performing a linear transformation, i.e. $w^T x + b$. We sample the weights then compute the output of the layer using the sampled weights and the mean of the biases.

`VariationalLogisticRegression` represents a logistic regression model using variational inference:

- The class is initialized with one linear variational layer used to perform the linear transformation in logistic regression.
- The `forward` method defines the forward pass for the logistic regression model by returning $f(x) = \sigma(w^T x + b)$ where σ is the sigmoid function.
- The `kl_divergence` method simply calls the same method of the `LinearVariational` layer to obtain the KL divergence term for the loss computation.

1.5. Comment the code of the training loop, especially the loss computation. Analyze the results provided by Figure 2.4. Compared to previous MAP estimate, how does the predictive distribution behave? What is the main difference between the Variational approximation and the Laplace approximation? The training loop uses a standard PyTorch format. The loss function calculates the Evidence Lower Bound (ELBO), which we want to maximize. In theory, we aim to maximize the likelihood of the data directly, but this is often intractable due to the integral over the weights. Therefore, we compute the Kullback-Leibler divergence $KL(q_{\theta}(w)||p(w))$ between the variational distribution $q_{\theta}(w)$ and the prior distribution $p(w)$. This acts as a regularization term, encouraging the variational distribution to be similar to the prior distribution. It represents the information lost when using $q_{\theta}(w)$ to approximate $p(w)$, which we want to minimize.

To ensure the model fits the data effectively, we compute the negative log-likelihood $NLL(\theta; \mathcal{D})$ of the data under the model parameterized by the weights sampled from $q_{\theta}(w)$. This is done using a binary cross-entropy loss. Subsequently, as maximizing ELBO is equivalent to minimizing:

$$\mathcal{L}_{VI}(\theta; \mathcal{D}) = NLL(\theta; \mathcal{D}) + KL(q_{\theta}(w)||p(w)),$$

we employ gradient descent to update the parameters of the variational distribution to better approximate the true posterior.

Compared to a MAP estimate, the variational approach does not just find the most probable weights (as MAP does) but instead approximates the entire posterior distribution over the weights. In the variational approach, the predictive distribution captures the model's uncertainty about its predictions. This approach results in boundary decisions between the MAP and Laplace approximation. The outcome of using the variational approach is a decision boundary that lies between what is obtained using MAP and the Laplace approximation. This results in a decision frontier that is distinct in the central area, yet it also effectively encompasses the aleatoric uncertainty. This approach provides a more comprehensive understanding of the model's predictions and uncertainties.

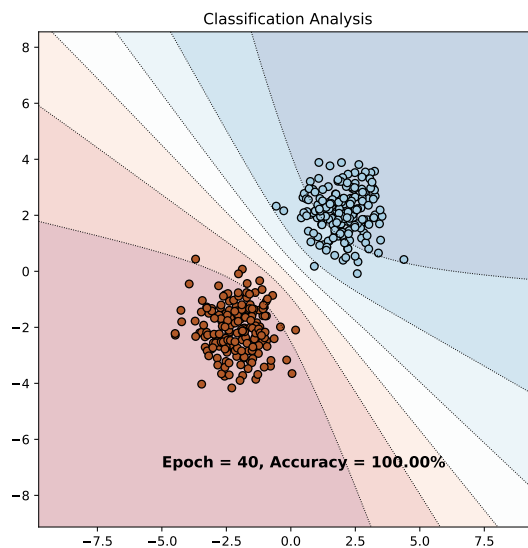


Figure 2.4: Illustration of a Variational Logistic Regression model applied to a binary classification task.

2.2 Bayesian Neural Networks

2.2.1 Variational Inference with Bayesian Neural Networks

2.1. Analyze the results showed on Figure 2.5. By applying Bayesian principles to a neural network with two hidden layers, we create a model capable of capturing intricate patterns within the data. In this context, each neuron's weight is treated as a random variable, representing our uncertainty about its true value. Consequently, we obtain a complex and non-linear decision boundary. Notably, the shaded regions, indicating the model's predictive uncertainty, reveal that the model demonstrates high confidence near the training data points and experiences increased uncertainty as it moves further away. Interestingly, this behavior leads to the emergence of "clusters" resembling the moon-shaped patterns found in the dataset. Moreover, this approach offers explainability, as it provides outputs that are not just binary but instead resemble the level of confidence one might have when predicting the locations of new data points.

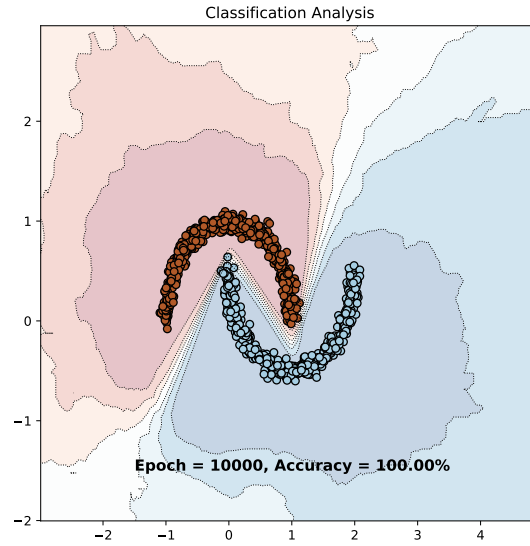


Figure 2.5: Illustration of a Bayesian Neural Network model applied to a binary classification task.

2.2.2 Monte Carlo Dropout

2.2. Again, analyze the results showed on Figure 2.6. What is the benefit of MC Dropout variational inference over Bayesian Logistic Regression with variational inference? Monte Carlo dropout is a technique that aligns with variational inference in Bayesian neural networks, where the dropout mechanism acts as a variational distribution for the network weights. Essentially, dropout introduces Bernoulli random variables, leading to a posterior predictive distribution that accounts for weight uncertainty. The formula for the predictive distribution of an output y for a new input x^* is:

$$p(y|x^*, \mathbf{X}, \mathbf{Y}) \approx \frac{1}{S} \sum_{s \in S} p(y^*|x^*, \mathbf{w}_s),$$

where \mathbf{w}_s represents the network weights after dropout, and S is the number of MC samples or dropout iterations.

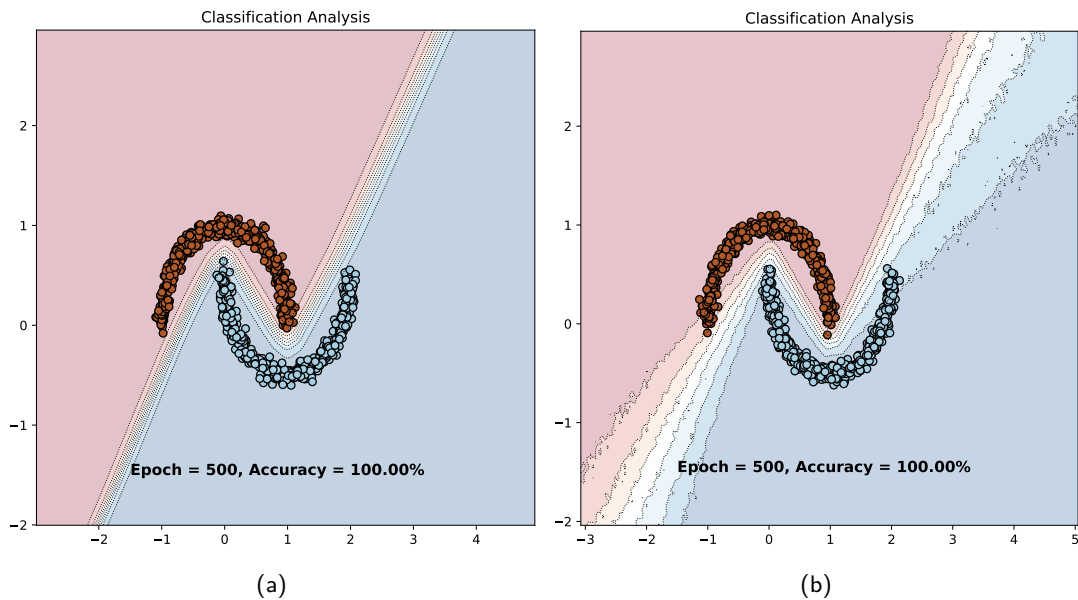


Figure 2.6: Illustration of a Bayesian Neural Network model applied to a binary classification task using (a) dropout and (b) Monte-Carlo dropout.

The results, shown in Figure 2.6b, include a decision boundary with adjacent bands indicating the model's confidence levels. These bands are formed by applying dropout during inference and averaging results from multiple stochastic passes. This approach illustrates the model's uncertainty in predictions, which contrast

to the smooth gradients of deterministic neural networks (shown in Figure 2.6a). The speckled appearance of uncertainty regions in the plot reflects how confidence varies across different input regions, due to the randomness introduced by MC Dropout. Incorporating MC Dropout in a standard neural network effectively turns it into Bayesian-like model, capable of expressing uncertainty in predictions. This approach not only allows the network to learn complex decision boundaries but also provides estimates of uncertainty, something often missing in regular neural networks. This probabilistic interpretation can lead to more informed decisions, as it shows how reliable the network's predictions are across different input areas.

Chapter 3

Uncertainty Applications

3.1 Monte-Carlo Dropout on MNIST

In this section, we focus on training a model using Monte Carlo Dropout (MC Dropout) on the MNIST dataset. Our primary aim was to identify the most uncertain samples. To do this, we used the variation ratio metric, which effectively measures epistemic uncertainty and is straightforward to calculate. For a given image, denoted as \mathbf{x} , we perform T stochastic forward passes through the model and record the predicted labels. We then determine the frequency $f_{\mathbf{x}}^{c^*}$ of the most common label (c^*) across the T passes. The variation ratio for image \mathbf{x} is calculated using the formula:

$$\text{var-ratio}[\mathbf{x}] = 1 - \frac{f_{\mathbf{x}}^{c^*}}{T}.$$

This formula provides a quantitative measure of uncertainty for an image.

1.1. What can you say about the images themselves? How do the histograms along them helps to explain failure cases? Finally, how do probabilities distribution of random images compare to the previous top uncertain images? In this experiment, we used a LeNet-5 styled model with Monte-Carlo dropout variational inference. The model was train on MNIST for 20 epoch using cross-entropy in a clascal way. Then we use the model to compute the variation ratios for each test image. This permit to retrieve images by them uncertainty. That what have been done in Figure 3.1, they denote five measurement over the models probabilities outputs for two certain and uncertain images. To sum up the behaviour of the outputed probabilities over the $T = 100$ stochastic forward pass, we used histograms that display four distribution. The first columns represent the distribution of the mean of the outputed probability per class. The second is the distribution of predicted class over the T forward passes. Then the 3 last columns represent the distribution of the outputed probability of a particular class (the most predicted class for the 3th column, the ground truth class in the 4th column, a another different class the the 5th column).

Those histogram represent how the output probabilities vary other $T = 100$ draw. If the model is not confident in its prediction, this will be reflected by high variation in outputed probabilities between each draw and so, histograms will be more sparce and the predicted choosen class will vary.

Figure 3.1 present images caracterized as certain by the model. With our humain eye, they seem to clearly represent their denoted class. About the presented distributions, they are composed of a single peak, this mean that we approximatly draw the same value every time. The mean probabilities for the predicted class is equal to one and all other class always have a probabilitiy of zero, meaning that the model is pretty confident about it's output.

In the other hand Figure 3.2 present the same thing for images caracterized as uncertain by the model. Dplayed numbers seem much more ambiguous, and event with our humain eye we have diffuculties to find out which number they represent. Distributions are much more spreaded. What does it mean for each columns? Does it mean the same for each columns? In the following part, we will try to explain this behaviour.

For the first two columns, this only mean that many different output class neurons can be activated. To explore more this behaviour we can look at the distribution of the predicted class. As we can see, both histograms are really similar. Cela interroge sur la valeur de la probabilité lorsque la classe est prédite.

The first interpretation that come in mind is that when a class neuron is activated, around his displayed mean value, it is almost the maximum and so the choosen class.

For the third and fourth column, this mean that they can furthermore take a wide range of values (large std). Those high variation in the output probabilities and in the choosen class translate an unconfident prediction.

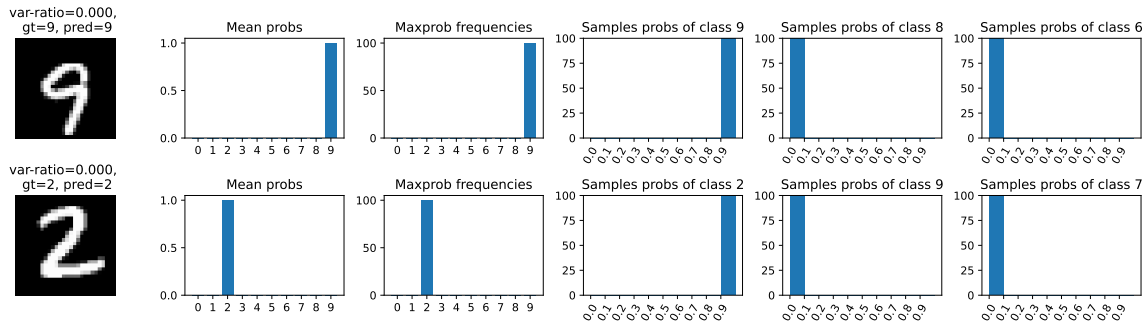


Figure 3.1

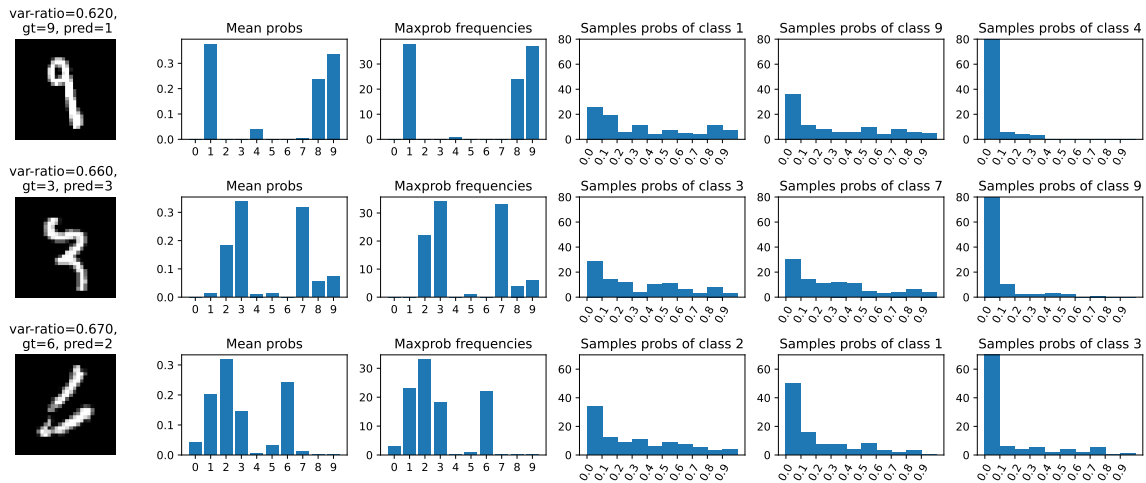


Figure 3.2

As we said when looking distribution of the mean of the probabilities of unconfident images,

3.2 Failure prediction

In this section, we have conducted a comparison of various methods aimed at obtaining a reliable confidence measure for model predictions. Such a measure permit to distinguish correct and incorrect prediction. An intelligent decision system equipped with such metrics can make informed choices, including adhering to the model's prediction or, conversely, involving a human operator, activating a backup system equipped with additional sensors, or triggering an alarm. This field of application is commonly referred to as failure prediction.

During the lecture, we found that Maximum Confidence Probability (MCP) is not a great metric for failure prediction. It assigns high confidence values to both correct and erroneous predictions because modern models tend to be overconfident, resulting in overlapping distributions between successes and errors. This issue persists even when using temperature scaling calibration.

Alternatively, when the model makes a misclassification, the probability associated with the true class y tends to be lower than the maximum probability, often falling to a low value. This observation leads us to consider the True Class Probability (TCP) as a suitable measure of uncertainty. However, the true class labels y are not available when estimating confidence for test inputs. This motivates the development of ConfidNet, whose primary objective is to directly regress the TCP value from the input image, allowing us to obtain a reliable measure of uncertainty without access to ground truth labels.

In this practical, we implemented ConfidNet to address failure predictions and compared it to two other methods that rely solely on the model's output probabilities. The first method is MCP, and the second is the entropy of the output probabilities. For these two methods, we used the previously trained MC Dropout model to compute the output probabilities. ConfidNet was trained for 30 epochs with the previous MC Dropout model as a teacher with frozen parameters and Mean Squared Error as the loss function.

2.1. Compare the precision-recall curves of each method along with their AUPR values. Why did we use AUPR metric instead of standard AUROC? To assess and compare these methods, we require

a suitable metric. Our objective is to identify classification errors, which we consider as the positive detection class, while correct predictions serve as the negative detection class. Since our models excel at making accurate predictions, we anticipate a low occurrence of classification errors, resulting in an imbalanced setting with a significant number of true negatives.

We have opted to employ the AUPR (Area Under the Precision-Recall Curve) instead of the AUROC (Area Under the Receiver Operating Characteristic Curve) due to the latter's unsuitability for imbalanced datasets. AUROC treats both classes equally and can yield misleading results, particularly under the influence of a large number of true negatives. This may overstate the model's performance, particularly in distinguishing the minority class, which in this context comprises classification errors. Conversely, AUPR is a more suitable choice, as it prioritizes precision and recall for the minority class. Precision measures the fraction of actual positives among the positive predictions, while recall measures the proportion of correctly identified actual positives. Consequently, AUPR proves to be a more reliable metric in our situation, where the positive class (classification errors) is significantly smaller than the negative class (correct predictions).

Figure 3.3 illustrates the precision-recall curves for each method, accompanied by their respective AUPR values. The results demonstrate that ConfidNet surpasses the other two methods. This superiority can be attributed to ConfidNet's training, which directly estimates the TCP value—a more dependable uncertainty measure compared to MCP and entropy, as elaborated upon in the preceding section. Notably, ConfidNet exhibits a slower decline in precision as recall increases, indicating a more balanced performance in terms of identifying true positives without a substantial rise in false positives.

For a better understanding of the lecture's confidence metrics, we aimed to compare the performance of predictive entropy and mutual information in the context of failure detection. The results are presented in Figure 3.4. It seems that entropy outperforms mutual information as a metric for failure detection. This distinction arises from the fact that predictive entropy measures aleatoric uncertainty, while mutual information assesses epistemic uncertainty. In the specific experiment of failure detection on MNIST, it proves more advantageous to rely on and measure aleatoric uncertainty, which originates from the natural variability or randomness of the numbers in the images. This explanation underscores why entropy is a superior metric to mutual information in this particular context.

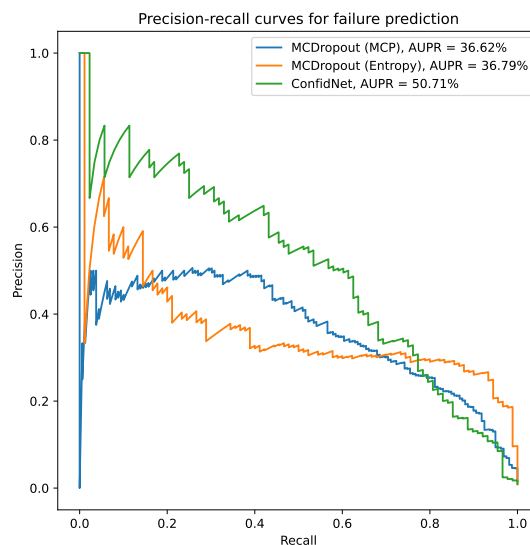


Figure 3.3

Figure 3.4

3.3 Out-of-distribution detection

Nowadays, models are trained on million of image (like imagenet), making epistemic uncertainty tend to zero. However, critical real word application like autonomous driving, medical imaging, etc. require to detect out-of-distribution (OOD) inputs because of the stochasticness of real life data. In this section, we will use Kuzushiji-MNIST (KMNIST) dataset as out of distribution data for our MC dropout MNIST predictor. This dataset is composed of 70,000 28x28 grayscale number images. We will keep precision, recall and AUPR as metrics to compare the different methods.

In particular in this section, we implemented the ODIN method (Liang et al., 2017) that enhance maximum softmax probabilities with temperature scaling and inverse adversarial perturbation. Those two technics are used to increase the difference between in and out of distribution.

The temperature scaling is a simple scaling of the logit by a temperature parameter $1/T$ before the softmax : $S_i(\mathbf{x}, T) = \frac{\exp(f_i(\mathbf{x})/T)}{\sum_{j=1}^N \exp(f_j(\mathbf{x}/T))}$ with S_i being the outputed probabilities, \mathbf{x} the input image, T the temperature parameter and f_i the logit of the i th class.

Inverse adversarial perturbation is to preprocess the input \mathbf{x} before feeding it to the neural network. The preprocessing is done by adding a small perturbation ϵ to the input image \mathbf{x} such that $\mathbf{x}' = \mathbf{x} - \epsilon \cdot \text{sign}(-\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, y))$ with \mathcal{L} the cross entropy loss function, y the ground truth label, ϵ the perturbation magnitude and \mathbf{x}' the perturbed image. The perturbation magnitude ϵ is chosen such that the perturbed image is still classified correctly by the neural network.

3.1. Compare the precision-recall curves of each OOD method along with their AUPR values. Which method perform best and why? Figure 3.5 illustrates the precision-recall curves for six uncertainty metrics (MCP, ODIN, CondifNet TCP, MC Dropout (MCD) mutual information and MCD predictive entropy), accompanied by their respective AUPR values. All precision-recall curves are smooth so we can mainly focus on AUPR values for our analysis.

ODIN with his two correction on output probabilities perform slightly better than MCP. It was expected as ODIN is a based on MCP. When comparing to the two last methods, ODIN is suprisingly outperformed by all them. This was not intended for mutual information and predictive entropy as they rely on the raw outputed probabilities as ODIN do. They seem to leverage the stochastic property of MC Dropout. In fact, those two last method performed really well with both 98% of AUPR, standing as the best metrics for OOD detection. Even with the leveraging of TCP, ConfidNet is behind with 96% of AUPR. Mutual information and predictive entropy with their simplicity seem to be the best metrics for OOD detection in our experiments but at the cost of having a heavier (bayesian) network. In fact in our case, computing entropy score take 18s whereas computing ODIN takes 4.1s.

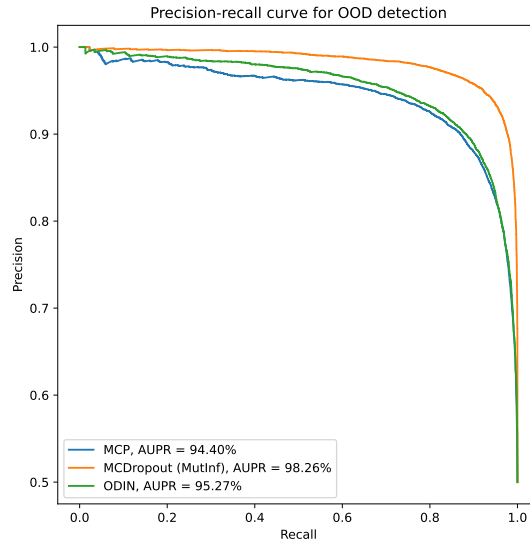


Figure 3.5

For further inversigation to observe the effect of ODIN technics, we tried to combine it with MC Dropout (by taking the mean of 100 stochastic forward passes), predictive entropy and mutual information. The results are presented in Figure 3.6. We can see that the combination of ODIN with MCD, predictive entropy and mutual information actually decrease performance of Entropy and mutual information, with a score 0.3% lower. Adding MCD to ODIN doesn't seem to have any effect.

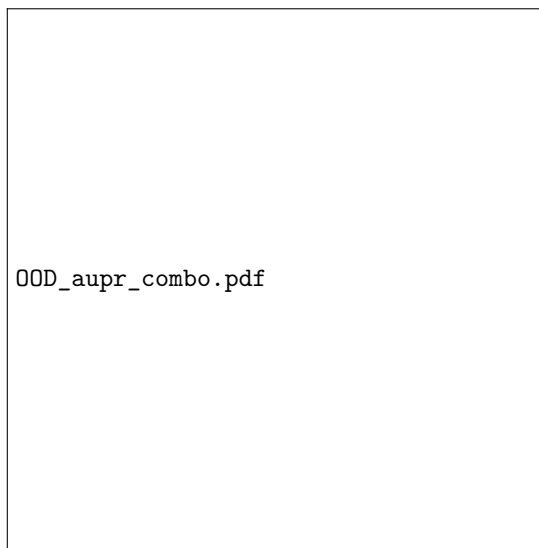


Figure 3.6

Seeing ODIN still failing, we thought it was because of wrong hyperparameter, in fact, the provided code used $\epsilon \in 0.006, 0.0006$ but suggested to use $\epsilon = 0.0014$ and the original paper sometime use high value of temperature, up to 10000. So we decided to conduct a grid search, AUPR value in fonction of ϵ and temperature are presented in Figure 3.7.

The overall effect of those parameters is good with a difference between the minimum and the maximum score of 1%. Setting the temperature to 10 give the best result with 98.42% of AUPR, beating MCD with entropy that had a score of 98.25%. Setting the temperature to 1 is equivalent to diable temperature scaling, this settings allow to find the effect of perturbation alone and vise versa for $\epsilon = 0$. Thus, the first cell represent a baseline for other cells. We can see by looking at the first line that temperature scaling allow to increase the AUPR by 1.1%, this is not that much but still a pretty good increase with our tight scores. We can apply the same logic for the effect of when $\epsilon = 0$. This time by looking at the first column, we cannot find any true effect of the perturbation with any value of ϵ . Thus, combining both parameter lead to the same score as if there was only temperature scaling. As the provided code use a smaller epsilon for the out of distribution data, we also tried the grid seach with the same scaling, it gave the same results.



Figure 3.7

Bibliography

Shiyu Liang, Yixuan Li, and R. Srikant. Principled detection of out-of-distribution examples in neural networks. *CoRR*, abs/1706.02690, 2017. URL <http://arxiv.org/abs/1706.02690>.