

# CS 6241 Final Project Report: Deep Kernel Learning for Satellite Imagery Classification

Aneesh Heintz

## 1 Introduction

Model prediction uncertainty is a crucial metric for various types of regression and classification tasks; especially in applications such as autonomous systems and medical diagnoses. Deep neural networks are flexible parametric models that can fit complex nonlinear patterns in data. Convolutional neural networks (CNN) are powerful pattern-recognition tools that have proven to be state of the art for image classification tasks. However, as they are increasingly implemented, their associated uncertainty estimates are not very accurate and do not produce confidence or uncertainty bounds on the predictions they create.

However, Bayesian inference has proven to be successful for learning under uncertainty. But the large number of parameters in CNNs means Bayesian inference is difficult to implement in the context of CNNs. Gaussian Processes (GP) are a powerful non-parametric tool in machine learning and allow predictions about data to be made by incorporating prior knowledge. GPs are easier suited for Bayesian inference. For a given training set of data, there are potentially infinitely many functions that fit the data. GPs assign probabilities to each function and the mean of the resulting probability distribution is the most probable characterization of the data. More importantly, GPs produce uncertainty estimates.

Using a deep neural network and Gaussian Process in sequence is known as deep kernel learning (DKL). DKL emerged as a useful research field to address the individual problems associated with neural networks and GPs [0]. It leverages the flexibility and interpretability uncertainty estimation framework of GPs with the ability to learn high-dimensional functions inherent in deep neural networks. Thus the neural network learns the kernel operator of a GP, which is in turn used to perform inference tasks [0].

This paper presents the creation and training of a convolutional deep kernel learning model to create classification predictions with associated uncertainty estimates. Therefore, the resulting model would reach performances similar to that of a CNN while adding a confidence metric of its prediction. Experiments are conducted on both synthetic and real datasets.

## 2 Related work

### 2.1 Neural Network Uncertainty

There are several methods that predate deep kernel learning that involve adding uncertainty bounds to predictions using GPs and CNNs. [1] shows that the output of a residual CNN with a 2D convolutional network prior over the weights and biases is a Gaussian Process in the limit of infinitely many convolutional filters. They show that state of the art and practical architectures such as CNNs and ResNets have equivalent GP representations. If each hidden layer has an infinite number of convolutional filters, the network prior is equivalent to a GP. [2] addresses uncertainty bounds of predictions using model dropout during training. They proposed a new theoretical framework that casts dropout training in the learning process of neural networks as Bayesian inference. This allows uncertainty estimates to be created directly from existing models.

### 2.2 Deep Kernel Learning

Deep Kernel learning was first proposed by [5] in an effort to combine the most useful assets associated with neural networks and GPs. GPs were popular because it is proven that bayesian neural nets with an infinite width converge to GPs with a particular kernel function. They are also flexible and interpretable. Neural networks, however, are superior in understanding representations in high-dimensional data. Therefore [5] trains a neural network and uses the top-level features of the deep neural network model as inputs to a GP. They show that DKL obtains results better than if its respective neural network model were trained on its own. The proposed DKL methods by [5] applied only to single-output regression problems and prohibits stochastic training. [6] expands upon [5] by proposing a new deep kernel learning model, Stochastic Variational DKL, that enables stochastic training, multi-task and multi-output learning. The SV-DKL architecture achieves competitive performance to stand-alone neural network architectures and beats the performance of similar methods that use neural networks and GPs in a coherent model.

### 3 Deep kernel learning method

The dataset being used contains  $n$  input images,  $X = x_1, \dots, x_n$ . Each image is described by a two-dimensional matrix with shape  $m \times m$ . Each image has an associated class category,  $y_i$ . The deep kernel learning applies an arbitrary convolutional neural network,  $h(x, w)$ , on the input images and is parameterized by weights  $w$ . The convolutional neural network creates embeddings,  $h_i(w, x)$  that are used as input features to a GP. For notation purposes, let the base kernel of a GP framework is defined as  $k(x_i, x_j | \theta)$ , where  $\theta$  is the parameters of the base kernel. Under the condition where the non-linear transformation,  $h(x, w)$  given by the neural network acts on the inputs of the entire model, the base kernel becomes  $k(w, h(x_i, w), h(x_j, w) | \theta)$ . The final layer of the convolutional neural network is passed as input to a Gaussian process:

$$f(h(x, w)) \sim \mathcal{GP}(\mu, k_\gamma)$$

where the mean vector,  $\mu_i = \mu(x_i)$ , and covariance matrix,  $(K_{h,h})_{ij} = k_\gamma(h(x_i, w), h(x_j, w))$ , is determined from the mean function and covariance kernel of the  $f(h(x, w))$ .

#### 3.1 Kernel (Covariance) Function

Then any collection of function values  $f$  has a joint Gaussian distribution,

$$f(h(x, w)) = [f(h_1), \dots, f(h_n)]^T \sim \mathcal{N}(\mu, K_{h,h})$$

The deep kernel learning model consists of a convolutional neural network followed by a Gaussian Process layer with as many outputs as there are classes in the dataset. The Gaussian Process layer uses an RBF kernel function:

$$k_{RBF}(x, x') = \text{cov}(f(x), f(x')) = a^2 \exp\left(-\frac{\|x - x'\|^2}{2l^2}\right)$$

where  $a$  and  $l$  are kernel hyper-parameters that controls the amplitudes and frequencies of the GP functions. The RBF kernel holds the assumption that function values at nearby inputs are more correlated than function values at far away inputs.

#### 3.2 Loss Function

The model is trained according to the following Variational Evidence Lower Bound (ELBO) loss function [7].

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{p_{\text{data}}(y, \mathbf{x})} [\mathbb{E}_{p(f|\mathbf{u}, \mathbf{x})q(\mathbf{u})} [\log p(y|f)]] - \beta \text{KL} [q(\mathbf{u}) \| p(\mathbf{u})]$$

Here,  $q(u)$  is the variational distributions for the inducing function values, and  $p(u)$  is the prior distribution for the inducing function values. The loss function was implemented as a maximum likelihood loss in gpytorch [7].

### 3.3 Training

The model is trained in an end-to-end manner optimized by SGD according to the variational Evidence Lower Bound loss. Learning rates and training duration varies based on the dataset the method is applied to. The model is written on the pytorch and gpytorch (developed at Cornell!) [7] frameworks and run on a Tesla K80 GPU using Google Colab.

## 4 Intermediate Goal: Cropped-region MNIST Dataset

### 4.1 Dataset relevance

The synthetic dataset used is the MNIST dataset [3] with regions of each image randomly cropped out. It was selected as an intermediate goal to demonstrate progress and that the overall method worked correctly in a "controlled environment." The dataset is provided through torchvision as separate training and validation sets. In total, 60 thousand training and 10 thousand validation samples are in each respective set.

### 4.2 Analysis

The architecture of the CNN used in the DKL model for training on the MNIST dataset is defined by a simple encoder. The DKL method is trained in an end-to-end manner and use a convolutional neural network in succession with a GP to obtain uncertainty estimates without decreasing performance. A GP uses the weights of the CNN's final layer, before the softmax activation function, as inputs. The DKL model creates a confidence interval for each class prediction. The model achieves a 94.5% accuracy when categorizing the dataset on previously unseen cropped handwriting samples.

Figure 1 demonstrates one sample prediction by the model. Figure 1(a) is the input handwriting sample with a randomly cropped region. It corresponds to a ground truth value of 8. Figure 1(b) shows the distribution of the DKL model's GP for the respective handwriting sample. The prediction output shows a prediction value along with the variance scores for each class.

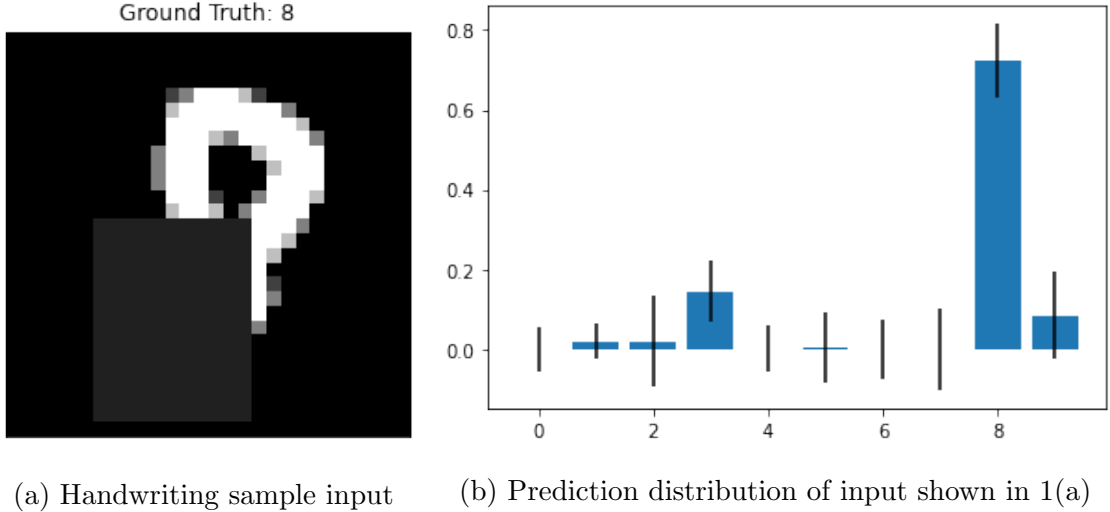


Figure 1: Sample prediction of the DKL model trained on cropped-region MNIST

## 5 UC Merced Land Use Dataset

The UC Merced Land Use dataset [4] is a collection of satellite images of different categories that appear in urban areas. The dataset includes images extracted from large images from the USGS National Map Urban Area Imagery collection from various urban areas around the United States. The dataset consists of 21 different classes, such as agriculture, beaches, harbors, runways, golf courses, etc. Each class is associated with 100 images of size  $256 \times 256$  pixels. Each pixel corresponds to 1 foot. Although this real-world dataset does not encompass all types of land regions that can be seen in satellite imagery, it nonetheless provides a good basis to demonstrate the proposed method.

### 5.1 Dataset relevance

Satellite imagery is useful for many reasons. Obviously, it is useful to gather information about the weather. Satellite imagery can aid tracking the spread of wildfires or the progression of hurricanes. Satellite imagery can also be useful for gathering economic information, such as measuring ship or airplane movements around the world. Measuring how busy each port or airport is can be quite informative and lucrative. This is especially important with regard to the current impacts associated with COVID-19. Satellite imagery has demonstrated the impact of the virus around the world. In Venice, for example, it can be seen through satellite imagery that the canals have cleared up and are far less polluted now that boats and ships aren't operating in the nearby waters. Lastly, satellite imagery is often used

for surveillance purposes. Both the government and human rights organizations use imagery of urban and rural areas to track vehicle movements and other activity at known locations over time.

The useful applications of satellite imagery are very broad. However, much of the aforementioned use cases for satellite imagery is analyzed manually. Automating classification, identification, and segmentation tasks on satellite imagery through methods like the one proposed to extract and analyze useful information could yield impressive results faster. When automated, such processes must have associated uncertainty estimates to be used in future decision making tasks.

## 5.2 Experiments

### 5.2.1 Performance comparison to a Stand-alone CNN

Using a simple CNN encoder, the DKL method achieves a performance accuracy of 75% on the UC Merced Land Use dataset. While the accuracy of the model is not as high as when it was applied to the synthetic MNIST dataset, it outperforms the same CNN encoder architecture trained end-to-end but without a GP. The comparative results are summarized in Table 1.

Table 1: Comparison between the DKL CNN-GP and a stand-alone CNN.

	Validation Loss	Validation Accuracy
<b>DKL CNN-GP</b>	1.18	75.20%
<b>DKL CNN-GP (noisy test data)</b>	1.28	70.99%
<b>Stand-alone CNN</b>	2.36	69.21%

An additional benefit of using the DKL CNN-GP over the stand-alone CNN is the confidence metrics associated with the predictions. For demonstration, Figure 2 shows a few selected example prediction distributions with respective confidence intervals from the DKL CNN-GP with a given satellite input image.

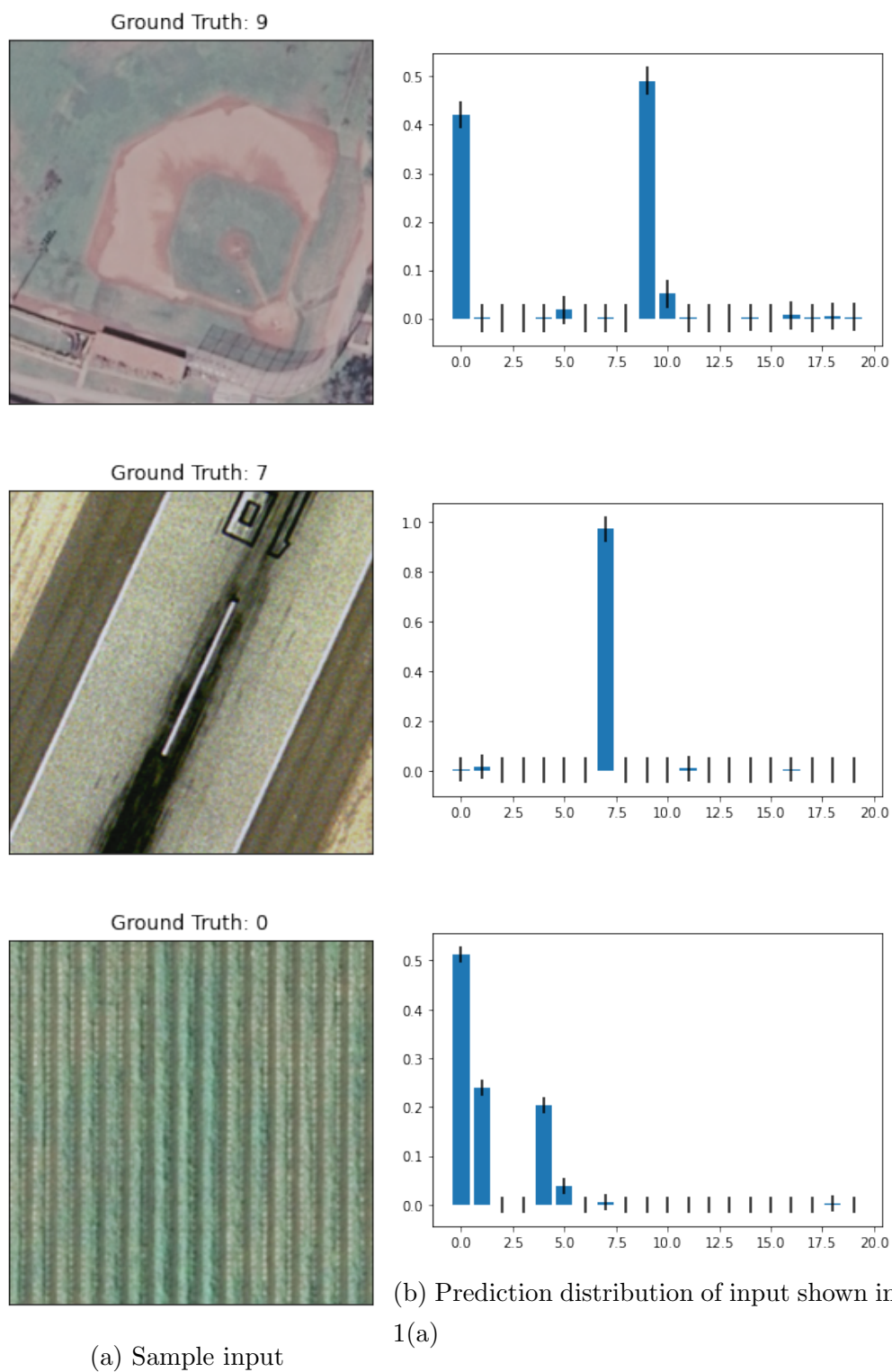


Figure 2: Sample predictions of the DKL model trained on the UC Merced Land Use dataset

### 5.2.2 Adding Noise

To demonstrate the usefulness of the confidence intervals associated with class predictions, a study was conducted on how uncertainty estimates varied with data that has a distribution that differs from the underlying distribution of the training data. As the UC Merced Land Use dataset has already been pre-processed to filter out noise, shot (Poisson) noise was added back in to each image in a validation set to simulate the effect that electronics on board the satellite would have on the camera during operation. This mimics an operational environment where the proposed method would be used on-board a spacecraft, something not currently done. On this dataset, the DKL method achieves a performance accuracy of 70.99%. Comparisons with the aforementioned experiments is detailed in Table 1. Selected prediction distributions are shown in Figure 3. They show an increase in the variance of the predictions, as expected when noise is added to images. However, it demonstrates that the proposed DKL CNN-GP method, when trained end-to-end, is able to achieve similarly competitive performances on a dataset that has a different underlying distribution from the training data.

### 5.2.3 Comparison to a Pre-trained CNN-GP

[6] demonstrated that training deep kernel learning methods end-to-end produced better results than training a GP with inputs derived from a pre-trained CNN. The last experiment performed further validated this conclusion. This experiment trained the DKL method in the same manner as the ones before. However, the CNN component of the DKL model was a ResNet18 module pre-trained on ImageNet. After training, the model achieved an accuracy of roughly 67%, similar to that of the stand-alone CNN.

## 6 Conclusion

In the robotics and space technology community, deep learning is a widely studied field. The impacts of implementing learning algorithms is wide. However, their implementations are not yet trustworthy, robust, and reliable due to the "black box" nature of neural networks and the lack of proper guidelines to properly and successfully train models. However, adding confidence metrics on deep learning predictions is a step in the right direction.

On both the synthetic MNIST dataset and the real-world Land Use dataset, the GP trained on the final layer of a simple encoder in an end-to-end fashion achieves a performance competitive to a simple stand-alone CNN. Additionally, this paper demonstrated that the end-to-end



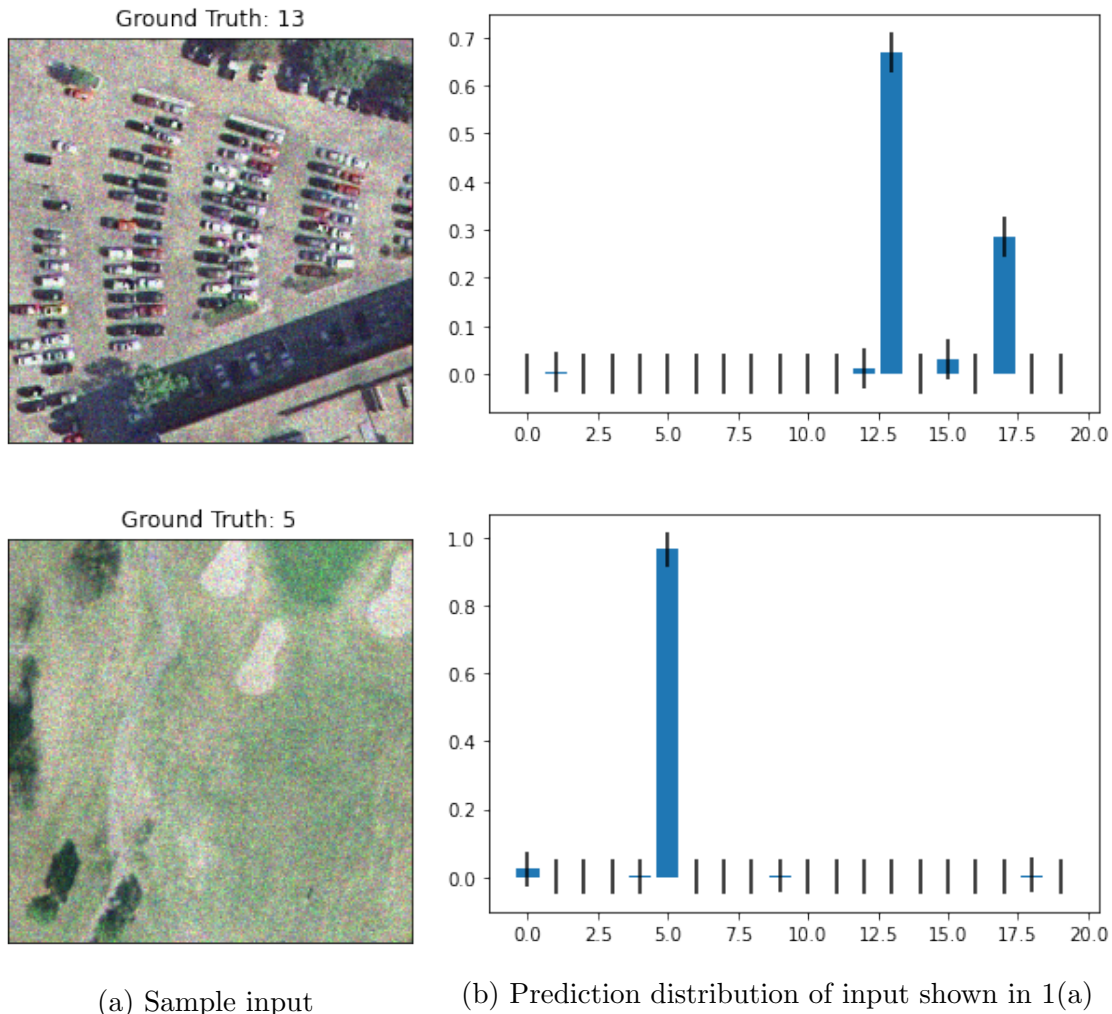


Figure 3: Results of shot noise added to the UC Merced Land Use dataset

training of the proposed DKL CNN-GP method is able to understand data with a different underlying distribution than what the model was trained on. The GP is also able to create confidence intervals associated with resulting class predictions from the neural network’s output weights. Confidence metrics are extremely powerful and informative. They can aid in a wide variety of tasks, such as autonomous decision making and streamlined analysis of data.

## 7 Project relevance / motivation

I am particularly interested in Gaussian Processes because it uses lazy learning to create a prediction that has uncertainty estimates associated with predictions.

My Ph.D. research area is in creating new ways for spacecraft to navigate relative to small objects using deep learning methods. While operating, the state estimates must be very precise and the measurements must be bounded by uncertainty estimates. Up until this course, I had not known that neural networks could have associated uncertainty estimates. I had a hard time creating deterministic uncertainty estimates for the models I have implemented so far. Therefore, I wanted to learn more about state of the art research that has been performed to model the correspondence between neural networks and Gaussian Processes. The addition of knowledge learned during this project will dramatically help the future direction and strength of the research I conduct.

Note: During the proposal stage of this project, I was not aware that deep kernel learning (combining neural networks followed by GPs with end-to-end training) was a research area. However, after starting the project and doing a more exhaustive literature review (see related work section), I amended my project to reference deep kernel learning.

## 8 References

- [0] Dasgupta, Sambarta & Sricharan, Kumar & Srivastava, Ashok. Finite Rank Deep Kernel Learning. NeurIPS, 2018. <http://bayesiandeeplearning.org/2018/papers/98.pdf>
- [1] Garriga-Alonso, Adrià & Aitchison, Laurence & Rasmussen, Carl. Deep Convolutional Networks as shallow Gaussian Processes. ICLR, 2019.
- [2] Gal, Yarin & Ghahramani, Zoubin. (2015). Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. Proceedings of The 33rd International Conference on Machine Learning.
- [3] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. "Gradient-based learning applied to document recognition." Proceedings of the IEEE, 86(11):2278-2324, November 1998.
- [4] Yi Yang and Shawn Newsam, "Bag-Of-Visual-Words and Spatial Extensions for Land-Use Classification," ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems (ACM GIS), 2010.
- [5] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. 2016. Deep

kernel learning. In Proceedings of the 19th International Conference on Artificial Intelligence and Statistics. 370–378.

[6] Andrew G Wilson, Zhiting Hu, Ruslan R Salakhutdinov, and Eric P Xing. Stochastic variational deep kernel learning. In Advances in Neural Information Processing Systems, pp. 2586–2594, 2016a.

[7] Gardner, Jacob R., Geoff Pleiss, David Bindel, Kilian Q. Weinberger, and Andrew Gordon Wilson. "GPYtorch: Blackbox Matrix-Matrix Gaussian Process Inference with GPU Acceleration." In Advances in Neural Information Processing Systems (2018)

## 9 Appendix: MNIST Code

```
1 ! pip install gpytorch
2 import torch
3 import torchvision
4 from torchvision import datasets, transforms
5 import matplotlib.pyplot as plt
6 import torch.nn as nn
7 import torch.nn.functional as F
8 from torch.optim import SGD, Adam
9 import torch.optim as optim
10 from torch.optim.lr_scheduler import StepLR
11 import tqdm
12 from torch.optim.lr_scheduler import MultiStepLR
13 import gpytorch
14 import math
15 import numpy as np
16
17 use_cuda = torch.cuda.is_available()
18 device = torch.device("cuda" if use_cuda else "cpu")
19 num_classes = 10
20 kwargs = {'num_workers': 1, 'pin_memory': True} if use_cuda else {}
21 args = {'log_interval': 1, 'learning_rate': 1e-1, 'step_size': 75, 'gamma': 0.1, 'n_epochs':
22         15, 'batch_size_train': 128, 'batch_size_test': 64}
23 torch.manual_seed(1)
24 print(use_cuda)
25
26 normalize = torchvision.transforms.Normalize((0.1307,), (0.3081,))
27 train_loader = torch.utils.data.DataLoader(
28     torchvision.datasets.MNIST('./files/', train=True, download=True,
29                               transform=torchvision.transforms.Compose([
30                                   torchvision.transforms.ToTensor(),
31                                   normalize,
32                                   transforms.RandomErasing()
33                               ])),
34     batch_size=args['batch_size_train'], shuffle=True)
```

```

34
35 test_loader = torch.utils.data.DataLoader(
36     torchvision.datasets.MNIST('/files/', train=False, download=True,
37                               transform=torchvision.transforms.Compose([
38                                   torchvision.transforms.ToTensor(),
39                                   normalize,
40                                   transforms.RandomErasing()
41                               ])),
42     batch_size=args['batch_size_test'], shuffle=True)
43
44 train_loader.dataset.train_data.shape, test_loader.dataset.train_data.shape
45
46 """Run CNN on MNIST"""
47
48 class CNN(nn.Module):
49     def __init__(self, num_features):
50         super(CNN, self).__init__()
51         self.conv1 = nn.Conv2d(1, 32, 3, 1)
52         self.conv2 = nn.Conv2d(32, 64, 3, 1)
53         self.dropout1 = nn.Dropout2d(0.25)
54         self.dropout2 = nn.Dropout2d(0.5)
55         self.fc1 = nn.Linear(9216, 128)
56         self.fc2 = nn.Linear(128, num_features)
57
58     def forward(self, x):
59         x = self.conv1(x)
60         x = F.relu(x)
61         x = self.conv2(x)
62         x = F.relu(x)
63         x = F.max_pool2d(x, 2)
64         x = self.dropout1(x)
65         x = torch.flatten(x, 1)
66         x = self.fc1(x)
67         x = F.relu(x)
68         x = self.dropout2(x)
69         x = self.fc2(x)
70         output = x
71         return output
72
73 class GaussianProcessLayer(gpytorch.models.ApproximateGP):
74     def __init__(self, num_dim, grid_bounds=(-10., 10.), grid_size=64):
75         variational_distribution = gpytorch.variational.CholeskyVariationalDistribution(
76             num_inducing_points=grid_size, batch_shape=torch.Size([num_dim])
77         )
78         variational_strategy = gpytorch.variational.MultitaskVariationalStrategy(
79             gpytorch.variational.GridInterpolationVariationalStrategy(
80                 self, grid_size=grid_size, grid_bounds=[grid_bounds],
81                 variational_distribution=variational_distribution,
82             ), num_tasks=num_dim,
83         )
84         super().__init__(variational_strategy)
85
86         self.covar_module = gpytorch.kernels.ScaleKernel(
87             gpytorch.kernels.RBFKernel(

```

```

88         lengthscale_prior=gpytorch.priors.SmoothedBoxPrior(
89             math.exp(-1), math.exp(1), sigma=0.1, transform=torch.exp
90         )
91     )
92 )
93 self.mean_module = gpytorch.means.ConstantMean()
94 self.grid_bounds = grid_bounds
95
96 def forward(self, x):
97     mean = self.mean_module(x)
98     covar = self.covar_module(x)
99     return gpytorch.distributions.MultivariateNormal(mean, covar)
100
101 class DKLModel(gpytorch.Module):
102     def __init__(self, feature_extractor, num_dim, likelihood, grid_bounds=(-10., 10.)):
103         super(DKLModel, self).__init__()
104         self.feature_extractor = feature_extractor
105         self.gp_layer = GaussianProcessLayer(num_dim=num_dim, grid_bounds=grid_bounds)
106         self.grid_bounds = grid_bounds
107         self.num_dim = num_dim
108         self.likelihood = likelihood
109
110     def forward(self, x):
111         features = self.feature_extractor(x)
112         features = gpytorch.utils.grid.scale_to_bounds(features, self.grid_bounds[0], self.
            grid_bounds[1])
113         features = features.transpose(-1, -2).unsqueeze(-1)
114         res_gp = self.gp_layer(features)
115         return res_gp
116
117 def train(args, model, likelihood, mll, device, train_loader, optimizer, epoch):
118     model.train()
119     likelihood.train()
120
121     total_loss = 0
122     loss_fn = nn.NLLLoss()
123
124     minibatch_iter = tqdm.notebook.tqdm(train_loader, desc=f"(Epoch {epoch}) Minibatch")
125     with gpytorch.settings.num_likelihood_samples(8):
126         for data, target in minibatch_iter:
127             data, target = data.to(device), target.to(device)
128             optimizer.zero_grad()
129             output = model(data)
130             loss = -mll(output, target)
131             total_loss += loss.item()
132             loss.backward()
133             optimizer.step()
134             minibatch_iter.set_postfix(loss=loss.item())
135     total_loss /= len(train_loader.dataset)
136     return total_loss
137
138 def test(model, likelihood, mll, device, test_loader):
139     model.eval()
140     likelihood.eval()

```

```

141     test_loss = 0
142     correct = 0
143     with torch.no_grad(), gpytorch.settings.num_likelihood_samples(16):
144         for data, target in test_loader:
145             data, target = data.to(device), target.to(device)
146             output = model(data)
147             loss = -mll(output, target)
148             test_loss += loss.item()
149             pred = likelihood(output).probs.mean(0).argmax(-1)
150             correct += pred.eq(target.view_as(pred)).sum().item()
151
152     test_loss /= len(test_loader.dataset)
153
154     return test_loss, correct / len(test_loader.dataset)
155
156 num_dim = 100
157 NN = CNN(num_dim)
158 likelihood = gpytorch.likelihoods.SoftmaxLikelihood(num_features=num_dim, num_classes=num_
    classes)
159 model = DKLModel(NN, num_dim, likelihood)
160
161
162 model = model.to(device)
163 likelihood = likelihood.to(device)
164 optimizer = SGD([
165     {'params': model.feature_extractor.parameters(), 'weight_decay': 1e-4},
166     {'params': model.gp_layer.hyperparameters(), 'lr': args['learning_rate'] * 0.01},
167     {'params': model.gp_layer.variational_parameters()},
168     {'params': model.likelihood.parameters()}],
169 lr=args['learning_rate'], momentum=0.9, nesterov=True, weight_decay=0)
170 scheduler = StepLR(optimizer, step_size=args['step_size'], gamma=args['gamma'])
171
172 mll = gpytorch.mlls.VariationalELBO(likelihood, model.gp_layer, num_data=len(train_loader))
173
174 for epoch in range(1, args['n_epochs'] + 1):
175     train_loss = train(args, model, likelihood, mll, device, train_loader, optimizer, epoch)
176     test_loss, acc = test(model, likelihood, mll, device, test_loader)
177     print('\n==> Epoch: {}, Train Loss: {:.4e}, Test Loss: {:.4e}, Test Acc: {:.4e}'.format(
178         epoch, train_loss, test_loss, acc))
179     scheduler.step()
180
181 """Individual Testing"""
182
183 examples = enumerate(test_loader)
184
185 batch_idx, (example_data, example_targets) = next(examples)
186 example_data.shape
187 idx = 8
188 fig = plt.figure()
189 plt.tight_layout()
190 plt.imshow(example_data[idx][0], cmap='gray', interpolation='none')
191 plt.title("Ground Truth: {}".format(example_targets[idx]))
192 plt.xticks([])
193 plt.yticks([])

```

```

193 plt.show()
194
195 model.eval()
196 data, target = example_data.to(device), example_targets.to(device)
197 output = model(data)
198 observed_pred = likelihood(output)
199 preds = observed_pred.probs.mean(0).cpu()
200 pred = preds.argmax(-1)
201 pred_distribution = preds.cpu().detach().numpy()
202 var = observed_pred.probs.var(1).cpu().detach().numpy()
203 pred.eq(example_targets.view_as(pred)).cpu().sum().item() / float(len(pred))
204
205 pred, example_targets
206
207 pred[idx].item(), example_targets[idx].item()
208
209 plt.figure()
210 plt.bar(np.arange(10), pred_distribution[idx].squeeze(), yerr=var[idx].squeeze())
211 plt.show()
212
213 pred_distribution[pred[idx].item()]
214
215 var[idx]
216
217 preds[0]

```

## 10 Appendix: Stand-alone CNN Code

```

1 ! pip install gpytorch
2 import torch
3 import torchvision
4 from torchvision import datasets, transforms
5 import torchvision.models as models
6 import matplotlib.pyplot as plt
7 import torch.nn as nn
8 import torch.nn.functional as F
9 from torch.optim import SGD, Adam
10 import torch.optim as optim
11 from torch.optim.lr_scheduler import StepLR
12 import tqdm
13 from torch.optim.lr_scheduler import MultiStepLR
14 import gpytorch
15 import math
16 import numpy as np
17 import os
18 from PIL import Image
19
20 # Commented out IPython magic to ensure Python compatibility.
21 from google.colab import drive
22 drive.mount('/content/drive/')
23 # % cd /content/drive/My Drive/Colab Notebooks/UCMerced_LandUse

```

```

24 use_cuda = torch.cuda.is_available()
25 device = torch.device("cuda" if use_cuda else "cpu")
26 kwargs = {'num_workers': 1, 'pin_memory': True} if use_cuda else {}
27 args = {'log_interval': 1, 'learning_rate': 1e-3, 'step_size': 50, 'gamma': 0.1, 'n_epochs':
        200, 'batch_size_train': 128, 'batch_size_test': 64}
28 torch.manual_seed(1)
29 print(use_cuda)
30
31 class Dataset(torch.utils.data.Dataset):
32     def __init__(self, mode, split):
33         path = './Images/'
34         categories = os.listdir(path)
35         imgs = []
36         target = []
37         for category_idx in (range(len(categories))):
38             img_categories = os.listdir(path + categories[category_idx])
39             for img_idx in (range(len(img_categories))):
40                 im_path = path + categories[category_idx] + '/' + img_categories[img_idx]
41                 imgs.append(im_path)
42                 target.append(category_idx)
43         self.imgs = np.array(imgs)
44         self.target = np.array(target)
45         self.img_exceptions = [130, 183, 209, 243, 396, 504, 505, 506, 507, 622, 623, 624, 633,
46                               770, 788, 858, 861, 863, 864, 865, 866, 867, 868, 869, 870, 915,
47                               935, 945, 993, 1055, 1060, 1077, 1122, 1145, 1146, 1308, 1320,
48                               1699, 1714, 1736, 1857, 2060, 2062, 2063]
49         self.img_idxxs = []
50         for idx in range(len(imgs)):
51             if idx not in self.img_exceptions:
52                 self.img_idxxs.append(idx)
53
54         self.img_idxxs = np.array(self.img_idxxs)
55         np.random.shuffle(self.img_idxxs)
56
57         self.imgs = self.imgs[self.img_idxxs]
58         self.target = self.target[self.img_idxxs]
59         self.num_classes = max(self.target)
60         self.categories = categories
61
62         split_idx = int(len(self.imgs) * split)
63         if mode == 'train':
64             self.imgs = self.imgs[:split_idx]
65             self.target = self.target[:split_idx]
66         elif mode == 'test':
67             self.imgs = self.imgs[split_idx:]
68             self.target = self.target[split_idx:]
69
70     def __len__(self):
71         return len(self.imgs)
72
73     def __getitem__(self, index):
74         idx = index
75
76         im_path = self.imgs[idx]

```



```

77     im = Image.open(im_path)
78     im = np.array(im) / 255.
79
80     means = [0.5, 0.5, 0.5]
81     stds = [0.5, 0.5, 0.5]
82
83     for ch in range(3):
84         im[:, :, ch] = (im[:, :, ch] - means[ch]) / stds[ch]
85
86     return im.reshape((3, 256, 256)), self.target[idx]
87
88 train_dataset = Dataset('train', 0.7)
89 test_dataset = Dataset('test', 0.7)
90 train_loader = torch.utils.data.DataLoader(train_dataset, shuffle = True)
91 test_loader = torch.utils.data.DataLoader(test_dataset, shuffle = True)
92
93 class CNN(nn.Module):
94     def __init__(self, num_features):
95         super(CNN, self).__init__()
96         self.conv1 = nn.Conv2d(3, 6, 5)
97         self.pool = nn.MaxPool2d(2, 2)
98         self.conv2 = nn.Conv2d(6, 16, 5)
99         self.fc1 = nn.Linear(59536, 120)
100        self.fc2 = nn.Linear(120, 84)
101        self.fc3 = nn.Linear(84, num_features)
102
103    def forward(self, x):
104        x = self.pool(F.relu(self.conv1(x)))
105        x = self.pool(F.relu(self.conv2(x)))
106        x = x.view(-1, 59536)
107        x = F.relu(self.fc1(x))
108        x = F.relu(self.fc2(x))
109        x = self.fc3(x)
110        x = F.log_softmax(x)
111        return x
112
113 class GaussianProcessLayer(gpytorch.models.ApproximateGP):
114     def __init__(self, num_dim, grid_bounds=(-10., 10.), grid_size=64):
115         variational_distribution = gpytorch.variational.CholeskyVariationalDistribution(
116             num_inducing_points=grid_size, batch_shape=torch.Size([num_dim])
117         )
118         variational_strategy = gpytorch.variational.MultitaskVariationalStrategy(
119             gpytorch.variational.GridInterpolationVariationalStrategy(
120                 self, grid_size=grid_size, grid_bounds=[grid_bounds],
121                 variational_distribution=variational_distribution,
122             ), num_tasks=num_dim,
123         )
124         super().__init__(variational_strategy)
125
126         self.covar_module = gpytorch.kernels.ScaleKernel(
127             gpytorch.kernels.RBFKernel(
128                 lengthscale_prior=gpytorch.priors.SmoothedBoxPrior(
129                     math.exp(-1), math.exp(1), sigma=0.1, transform=torch.exp
130                 )

```

```

131         )
132     )
133     self.mean_module = gpytorch.means.ConstantMean()
134     self.grid_bounds = grid_bounds
135
136     def forward(self, x):
137         mean = self.mean_module(x)
138         covar = self.covar_module(x)
139         return gpytorch.distributions.MultivariateNormal(mean, covar)
140
141 class DKLModel(gpytorch.Module):
142     def __init__(self, feature_extractor, num_dim, likelihood, grid_bounds=(-10., 10.)):
143         super(DKLModel, self).__init__()
144         self.feature_extractor = feature_extractor
145         self.gp_layer = GaussianProcessLayer(num_dim=num_dim, grid_bounds=grid_bounds)
146         self.grid_bounds = grid_bounds
147         self.num_dim = num_dim
148         self.likelihood = likelihood
149         self.drop = nn.Dropout(0.5)
150
151     def forward(self, x):
152         features = self.feature_extractor(x)
153         features = gpytorch.utils.grid.scale_to_bounds(features, self.grid_bounds[0], self.
            grid_bounds[1])
154         features = features.transpose(-1, -2).unsqueeze(-1)
155         res_gp = self.gp_layer(features)
156         return res_gp
157
158 def train(args, model, likelihood, mll, device, train_loader, optimizer, epoch):
159     model.train()
160
161     total_loss = 0
162     loss_fn = nn.NLLLoss()
163
164     minibatch_iter = tqdm.notebook.tqdm(train_loader, desc=f"(Epoch {epoch}) Minibatch")
165     with gpytorch.settings.num_likelihood_samples(8):
166         for data, target in minibatch_iter:
167             data = data.to(device).float()
168             target = target.to(device)
169             optimizer.zero_grad()
170             output = model(data)
171             loss = F.nll_loss(output, target)
172             total_loss += loss.item()
173             loss.backward()
174             optimizer.step()
175             minibatch_iter.set_postfix(loss=loss.item())
176     total_loss /= len(train_loader.dataset)
177     return total_loss
178
179 def test(model, likelihood, mll, device, test_loader, epoch):
180     model.eval()
181     #likelihood.eval()
182     test_loss = 0
183     correct = 0

```

```

184 minibatch_iter = tqdm.notebook.tqdm(test_loader, desc=f"(Epoch {epoch}) Minibatch")
185 with torch.no_grad(), gpytorch.settings.num_likelihoood_samples(16):
186     for data, target in minibatch_iter:
187         data = data.to(device).float()
188         target = target.to(device)
189         output = model(data)
190         loss = F.nll_loss(output, target)
191         test_loss += loss.item()
192         pred = output.argmax(dim=1, keepdim=True)
193         correct += pred.eq(target.view_as(pred)).sum().item()
194
195     test_loss /= len(test_loader.dataset)
196
197     return test_loss, correct / len(test_loader.dataset)
198
199 num_dim = 10
200 num_classes = train_loader.dataset.num_classes
201 model = CNN(num_classes)
202
203 model = model.to(device)
204 optimizer = SGD(model.parameters(), lr=args['learning_rate'], momentum=0.9, nesterov=True,
205                 weight_decay=0)
206 scheduler = StepLR(optimizer, step_size=args['step_size'], gamma=args['gamma'])
207 print(num_classes)
208
209 for epoch in range(1, args['n_epochs'] + 1):
210     train_loss = train(args, model, None, None, device, train_loader, optimizer, epoch)
211     test_loss, acc = test(model, None, None, device, test_loader, epoch)
212     print('\n==> Epoch: {}, Train Loss: {:.4e}, Test Loss: {:.4e}, Test Acc: {:.4e}'.format(
213         epoch, train_loss, test_loss, acc))
214     scheduler.step()
215
216 """Individual Testing"""
217
218 model.eval()
219 examples = iter(test_loader)
220 example_data, example_targets = next(examples)
221 data, target = example_data.to(device).float().view(-1,3,256,256), example_targets.to(device)
222
223 output = model(data)
224 observed_pred = likelihood(output)
225 preds = observed_pred.probs.mean(0).cpu()
226 pred = preds.argmax(-1)
227 pred_distribution = preds.cpu().detach().numpy()
228 var = observed_pred.probs.view((-1,train_dataset.num_classes)).var(1).cpu().detach().numpy()
229 print(pred.eq(example_targets.view_as(pred)).cpu().sum().item() / float(len(pred)))
230
231 idx = 0
232 fig = plt.figure()
233 plt.tight_layout()
234 im = example_data[idx].view((256,256,-1)) * 0.5 + 0.5
235 plt.imshow(im)
236 plt.title("Ground Truth: {}".format(example_targets[idx]))
237 plt.xticks([])

```

```

235 plt.yticks([])
236 plt.show()
237
238 plt.figure()
239 plt.bar(np.arange(train_dataset.num_classes), pred_distribution[idx].squeeze(), yerr=var[idx]
        ].squeeze())
240 plt.show()
241
242 pred[idx].item(), example_targets[idx].item()
243
244 train_dataset.categories[pred[idx].item()], train_dataset.categories[example_targets[idx].
        item()]
245
246 pred_distribution[idx]
247
248 var

```

## 11 Appendix: DKL CNN-GP Code

```

1 ! pip install gpytorch
2 import torch
3 import torchvision
4 from torchvision import datasets, transforms
5 import torchvision.models as models
6 import matplotlib.pyplot as plt
7 import torch.nn as nn
8 import torch.nn.functional as F
9 from torch.optim import SGD, Adam
10 import torch.optim as optim
11 from torch.optim.lr_scheduler import StepLR
12 import tqdm
13 from torch.optim.lr_scheduler import MultiStepLR
14 import gpytorch
15 import math
16 import numpy as np
17 import os
18 from PIL import Image
19
20 # Commented out IPython magic to ensure Python compatibility.
21 from google.colab import drive
22 drive.mount('/content/drive/')
23 # % cd /content/drive/My Drive/Colab Notebooks/UCMerced_LandUse
24 use_cuda = torch.cuda.is_available()
25 device = torch.device("cuda" if use_cuda else "cpu")
26 kwargs = {'num_workers': 1, 'pin_memory': True} if use_cuda else {}
27 args = {'log_interval': 1, 'learning_rate': 1e-3, 'step_size': 75, 'gamma': 0.1, 'n_epochs':
        50, 'batch_size_train': 128, 'batch_size_test': 64}
28 torch.manual_seed(1)
29 print(use_cuda)
30
31 class Dataset(torch.utils.data.Dataset):

```

```

32 def __init__(self, mode, split, noise=False):
33     path = './Images/'
34     self.noise = noise
35     categories = os.listdir(path)
36     imgs = []
37     target = []
38     for category_idx in (range(len(categories))):
39         img_categories = os.listdir(path + categories[category_idx])
40         for img_idx in (range(len(img_categories))):
41             im_path = path + categories[category_idx] + '/' + img_categories[img_idx]
42             imgs.append(im_path)
43             target.append(category_idx)
44     self.imgs = np.array(imgs)
45     self.target = np.array(target)
46     self.img_exceptions = [130, 183, 209, 243, 396, 504, 505, 506, 507, 622, 623, 624, 633,
47                             770, 788, 858, 861, 863, 864, 865, 866, 867, 868, 869, 870, 915,
48                             935, 945, 993, 1055, 1060, 1077, 1122, 1145, 1146, 1308, 1320,
49                             1699, 1714, 1736, 1857, 2060, 2062, 2063]
50     self.img_idxxs = []
51     for idx in range(len(imgs)):
52         if idx not in self.img_exceptions:
53             self.img_idxxs.append(idx)
54
55     self.img_idxxs = np.array(self.img_idxxs)
56     np.random.shuffle(self.img_idxxs)
57
58     self.imgs = self.imgs[self.img_idxxs]
59     self.target = self.target[self.img_idxxs]
60     self.num_classes = max(self.target)
61     self.categories = categories
62
63     split_idx = int(len(self.imgs) * split)
64     if mode == 'train':
65         self.imgs = self.imgs[:split_idx]
66         self.target = self.target[:split_idx]
67     elif mode == 'test':
68         self.imgs = self.imgs[split_idx:]
69         self.target = self.target[split_idx:]
70
71 def __len__(self):
72     return len(self.imgs)
73
74 def __getitem__(self, index):
75     idx = index
76
77     im_path = self.imgs[idx]
78     im = Image.open(im_path)
79     im = np.array(im)
80     if self.noise:
81         pois = 50
82         im = np.random.poisson(im / 255. * pois) / pois * 255
83
84     im = im / 255.
85

```

```

86     means = [0.485, 0.456, 0.406]
87     stds = [0.229, 0.224, 0.225]
88
89     means = [0.5, 0.5, 0.5]
90     stds = [0.5, 0.5, 0.5]
91
92     for ch in range(3):
93         im[:, :, ch] = (im[:, :, ch] - means[ch]) / stds[ch]
94
95     return im.reshape((3, 256, 256)), self.target[idx]
96
97 train_dataset = Dataset('train', 0.7)
98 test_dataset = Dataset('test', 0.7)
99
100 train_loader = torch.utils.data.DataLoader(train_dataset, shuffle = True)
101 test_loader = torch.utils.data.DataLoader(test_dataset, shuffle = True)
102
103 test_dataset_noise = Dataset('test', 0.7, noise = True)
104 test_loader_noise = torch.utils.data.DataLoader(test_dataset_noise, shuffle = True)
105
106 class CNN(nn.Module):
107     def __init__(self, num_features):
108         super(CNN, self).__init__()
109         self.conv1 = nn.Conv2d(3, 6, 5)
110         self.pool = nn.MaxPool2d(2, 2)
111         self.conv2 = nn.Conv2d(6, 16, 5)
112         self.fc1 = nn.Linear(59536, 120)
113         self.fc2 = nn.Linear(120, 84)
114         self.fc3 = nn.Linear(84, num_features)
115
116     def forward(self, x):
117         x = self.pool(F.relu(self.conv1(x)))
118         x = self.pool(F.relu(self.conv2(x)))
119         x = x.view(-1, 59536)
120         x = F.relu(self.fc1(x))
121         x = F.relu(self.fc2(x))
122         x = self.fc3(x)
123         return x
124
125 class GaussianProcessLayer(gpytorch.models.ApproximateGP):
126     def __init__(self, num_dim, grid_bounds=(-10., 10.), grid_size=64):
127         variational_distribution = gpytorch.variational.CholeskyVariationalDistribution(
128             num_inducing_points=grid_size, batch_shape=torch.Size([num_dim])
129         )
130         variational_strategy = gpytorch.variational.MultitaskVariationalStrategy(
131             gpytorch.variational.GridInterpolationVariationalStrategy(
132                 self, grid_size=grid_size, grid_bounds=[grid_bounds],
133                 variational_distribution=variational_distribution,
134             ), num_tasks=num_dim,
135         )
136         super().__init__(variational_strategy)
137
138         self.covar_module = gpytorch.kernels.ScaleKernel(
139             gpytorch.kernels.RBFKernel(

```

```

140         lengthscale_prior=gpytorch.priors.SmoothedBoxPrior(
141             math.exp(-1), math.exp(1), sigma=0.1, transform=torch.exp
142         )
143     )
144 )
145 self.mean_module = gpytorch.means.ConstantMean()
146 self.grid_bounds = grid_bounds
147
148 def forward(self, x):
149     mean = self.mean_module(x)
150     covar = self.covar_module(x)
151     return gpytorch.distributions.MultivariateNormal(mean, covar)
152
153 class DKLModel(gpytorch.Module):
154     def __init__(self, feature_extractor, num_dim, likelihood, grid_bounds=(-10., 10.)):
155         super(DKLModel, self).__init__()
156         self.feature_extractor = feature_extractor
157         self.gp_layer = GaussianProcessLayer(num_dim=num_dim, grid_bounds=grid_bounds)
158         self.grid_bounds = grid_bounds
159         self.num_dim = num_dim
160         self.likelihood = likelihood
161         self.drop = nn.Dropout(0.5)
162
163     def forward(self, x):
164         features = self.feature_extractor(x)
165         features = gpytorch.utils.grid.scale_to_bounds(features, self.grid_bounds[0], self.
            grid_bounds[1])
166         features = features.transpose(-1, -2).unsqueeze(-1)
167         res_gp = self.gp_layer(features)
168         return res_gp
169
170 def train(args, model, likelihood, mll, device, train_loader, optimizer, epoch):
171     model.train()
172     likelihood.train()
173
174     total_loss = 0
175     loss_fn = nn.NLLLoss()
176
177     minibatch_iter = tqdm.notebook.tqdm(train_loader, desc=f"(Epoch {epoch}) Minibatch")
178     with gpytorch.settings.num_likelihood_samples(8):
179         for data, target in minibatch_iter:
180             data = data.to(device).float()
181             target = target.to(device)
182             optimizer.zero_grad()
183             output = model(data)
184             loss = -mll(output, target)
185             total_loss += loss.item()
186             loss.backward()
187             optimizer.step()
188             minibatch_iter.set_postfix(loss=loss.item())
189     total_loss /= len(train_loader.dataset)
190     return total_loss
191
192 def test(model, likelihood, mll, device, test_loader, epoch):

```

```

193     model.eval()
194     likelihood.eval()
195     test_loss = 0
196     correct = 0
197     minibatch_iter = tqdm.notebook.tqdm(test_loader, desc=f"(Epoch {epoch}) Minibatch")
198     with torch.no_grad(), gpytorch.settings.num_likelihoood_samples(16):
199         for data, target in minibatch_iter:
200             data = data.to(device).float()
201             target = target.to(device)
202             output = model(data)
203             loss = -mll(output, target)
204             test_loss += loss.item()
205             pred = likelihood(output).probs.mean(0).argmax(-1)
206             correct += pred.eq(target.view_as(pred)).sum().item()
207
208     test_loss /= len(test_loader.dataset)
209
210     return test_loss, correct / len(test_loader.dataset)
211
212 num_dim = 100
213 num_classes = train_loader.dataset.num_classes
214 #resnet18 = models.resnet18(pretrained=True)
215 #NN = resnet18
216 NN = CNN(num_dim)
217 likelihood = gpytorch.likelihoods.SoftmaxLikelihood(num_features=num_dim, num_classes=num_
    classes)
218 model = DKLModel(NN, num_dim, likelihood)
219
220 model = model.to(device)
221 likelihood = likelihood.to(device)
222
223 learnable_params = [
224     {'params': model.feature_extractor.parameters(), 'weight_decay': 1e-6},
225     {'params': model.gp_layer.hyperparameters(), 'lr': args['learning_rate'] * 0.01},
226     {'params': model.gp_layer.variational_parameters()},
227     {'params': model.likelihood.parameters()}],
228 ]
229
230 optimizer = SGD(learnable_params, lr=args['learning_rate'], momentum=0.9, nesterov=True,
    weight_decay=0)
231 scheduler = StepLR(optimizer, step_size=args['step_size'], gamma=args['gamma'])
232
233 mll = gpytorch.mlls.PredictiveLogLikelihood(likelihood, model.gp_layer, len(train_loader))
234 print(num_classes)
235
236 for epoch in range(1, args['n_epochs'] + 1):
237     train_loss = train(args, model, likelihood, mll, device, train_loader, optimizer, epoch)
238     test_loss, acc = test(model, likelihood, mll, device, test_loader, epoch)
239     test_loss_noise, acc_noise = test(model, likelihood, mll, device, test_loader_noise, epoch
    )
240     print('\n==> Epoch: {}, Train Loss: {:.4e}, Test Loss: {:.4e}, Test Acc: {:.4e}, Test Loss
    Noise: {:.4e}, Test Acc Noise: {:.4e}'.format(epoch, train_loss, test_loss, acc, test
    _loss_noise, acc_noise))
241     scheduler.step()

```



```

242
243 """Individual Testing"""
244
245 model.eval()
246 examples = iter(test_loader_noise)
247 example_data, example_targets = next(examples)
248 data, target = example_data.to(device).float().view(-1,3,256,256), example_targets.to(device
    )
249 output = model(data)
250 observed_pred = likelihood(output)
251 preds = observed_pred.probs.mean(0).cpu()
252 pred = preds.argmax(-1)
253 pred_distribution = preds.cpu().detach().numpy()
254 var = observed_pred.probs.view((-1,train_dataset.num_classes)).var(1).cpu().detach().numpy()
255 print(pred.eq(example_targets.view_as(pred)).cpu().sum().item() / float(len(pred)))
256
257 idx = 0
258 fig = plt.figure()
259 plt.tight_layout()
260 im = example_data[idx].view((256,256,-1)) * 0.5 + 0.5
261 plt.imshow(im)
262 plt.title("Ground Truth: {}".format(example_targets[idx]))
263 plt.xticks([])
264 plt.yticks([])
265 plt.show()
266
267 plt.figure()
268 plt.bar(np.arange(train_dataset.num_classes), pred_distribution[idx].squeeze(), yerr=var[idx
    ].squeeze())
269 plt.show()
270
271 pred[idx].item(), example_targets[idx].item()
272
273 train_dataset.categories[pred[idx].item()], train_dataset.categories[example_targets[idx].
    item()]
274
275 pred_distribution[idx]
276
277 var

```