

Goals and Challenges

Long-Term Goal: To discover fundamental principles of learning, and design AI systems that continue to learn and improve throughout their life (like humans)

Living beings can learn throughout their life from small chunks of data in a non-stationary world, but deep learning requires a large amount of data from a stationary world. Our current research focuses on reducing this gap.

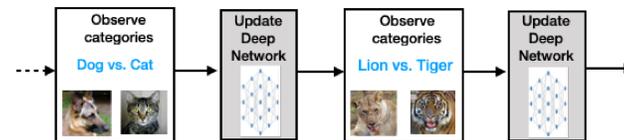
Main Idea: We use Bayesian principles to enable human-like learning of deep networks.

- New methods for uncertainty estimation in deep networks
- Convert deep networks to Gaussian process to use them as prior for life-long learning.

Standard Deep Learning

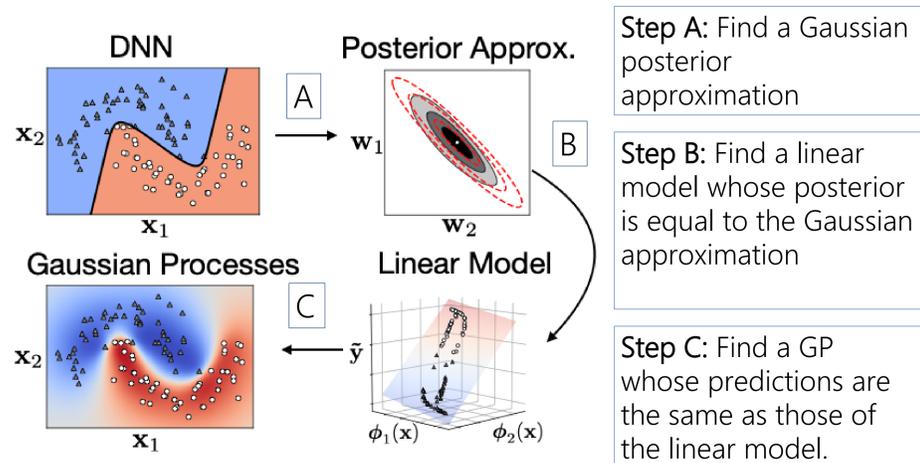


Deep Continual Learning

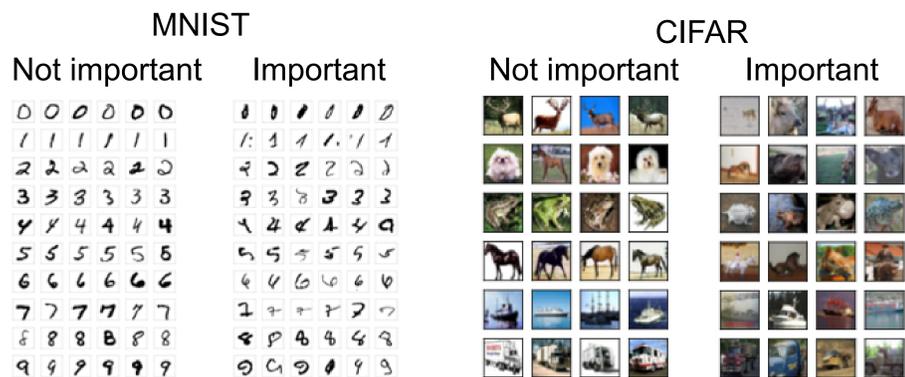


From Deep Neural Networks to Gaussian Processes

Using Gaussian posterior approximations, we can convert Neural network to Gaussian processes (GP). This enables us to use deep networks as functional priors.



Using this approach, we can identify important data examples.

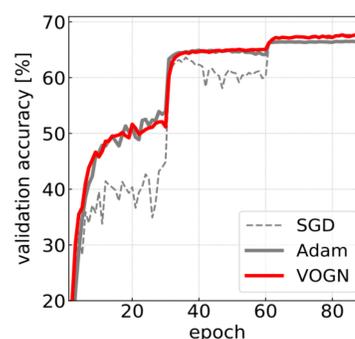


Approximate Inference Turns Deep Networks into Gaussian Processes. Khan ME, Immer A, Abedi E, Korzepa M (NeurIPS 2019).

Uncertainty Estimation for Large Deep Networks

- A new scalable algorithm, called **VOGN**, to estimate uncertainty for large deep-learning problems.
- For the first time, we can train ResNet-18 on ImageNet with 128 GPUs. We achieve similar performance to Adam/SGD in about the same number of epochs while preserving the benefits of Bayesian principles:

- predictive probabilities are well-calibrated
- uncertainties on out-of-distribution data are improved
- continual-learning performance is boosted



```

Algorithm Distributed VOGN
1: repeat
2:   Sample a minibatch M of size M.
3:   Split M into each GPU (local minibatch M_local).
4:   for each GPU in parallel do
5:     for k = 1, 2, ..., K (# MC samples) do
6:       Sample weight w^(k) ~ q(w).
7:       Compute gradient g^(k) w.r.t. w.
8:       g_k ← (1/M) * sum_{i in M_local} g_i^(k) and h_k ← (1/M) * sum_{i in M_local} (g_i^(k))^2.
9:     end for
10:    g ← (1/K) * sum_{k=1}^K g_k and h ← (1/K) * sum_{k=1}^K h_k.
11:  end for
12:  AllReduce (aggregate) g, h among all GPUs.
13:  m ← beta_1 * m + (1 - beta_1) * g and s ← (1 - tau * beta_2) * s + tau * h.
14:  mu ← mu - alpha * m / (s + delta).
15: until stopping criterion is met = 0
    
```

Bayesian Uncertainty Estimation in Image Segmentation



Segmentation (left) and its predictive entropy (right) on Cityscapes dataset

Practical Deep Learning with Bayesian Principles. Osawa K, Swaroop S, Jain A, Eschenhagen R, Turner RE, Yokota R, Khan ME (NeurIPS 2019).

A New Bayesian Learning Rule for Mixture of Exponential Family

- Bayesian Learning rule is a general learning rule from which many learning-algorithms can be derived
- Deep learning (SGD and Adam)
- Least-squares, Kalman filters etc.
- And many more..
- We extended the application of rule to mixture of exponential family, enabling us to derive, e.g., an ensemble Newton method.

Fast and Simple Natural-Gradient Variational Inference with Mixture of Exponential-family Approximations. Lin W, Khan ME, Schmidt M (ICML 2019).

Inference Networks for Gaussian-Process Models

- This work uses neural networks to approximate the posterior distribution of Gaussian process models.
- The main idea is to use a "functional" mirror descent algorithm.
- The updates can be approximated using a neural networks as the posterior distribution.

Scalable Training of Inference Networks for Gaussian-Process Models. Shi J, Khan ME, Zhu J (ICML 2019).

A Generalization Bound for Online Variational Inference

- This work derives new bounds for methods that perform online variational inference.
- We considers a variety of such algorithms and showed that the generalization error is bounded, even though each step of the algorithm is approximate.
- **Best paper award at ACML 2019.**

A Generalization Bound For Online Variational Inference. Chérif-Abdellatif, BE, Alquier, P and Khan, ME (ACML 2019).