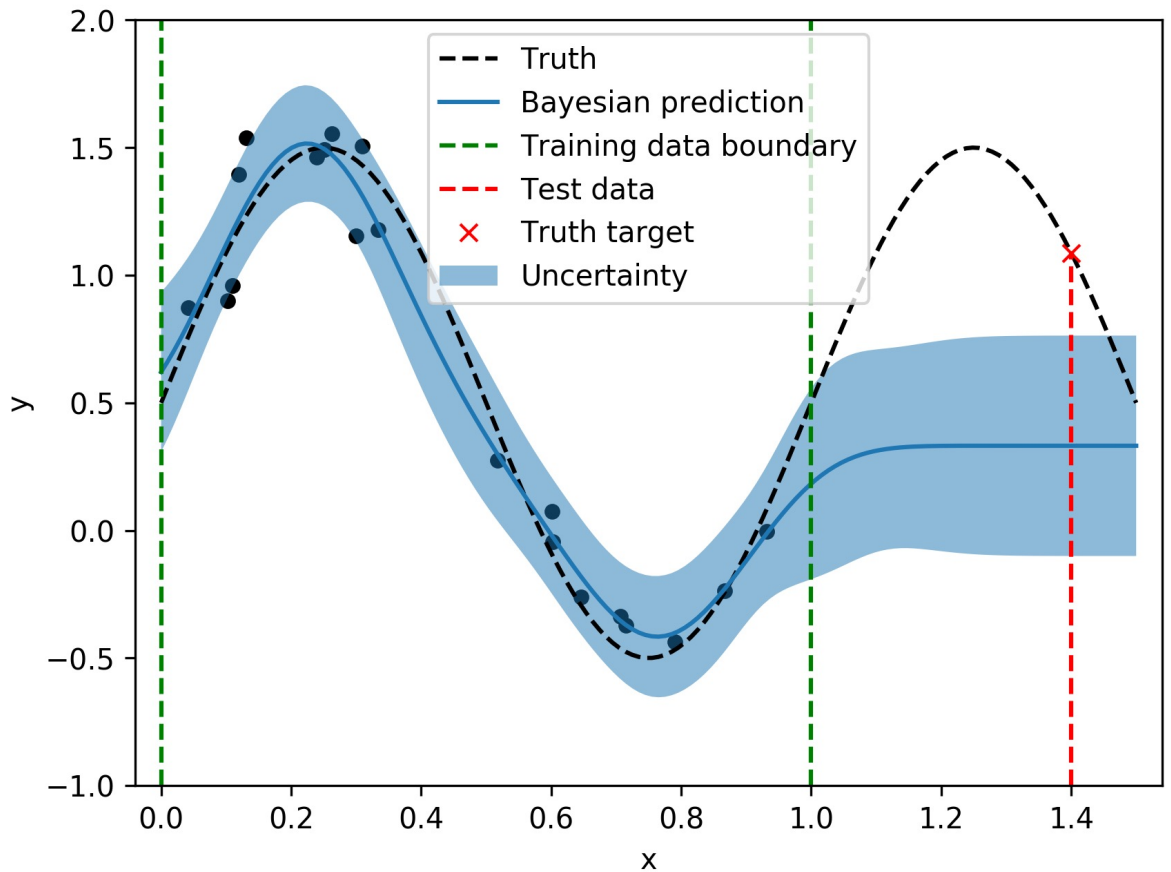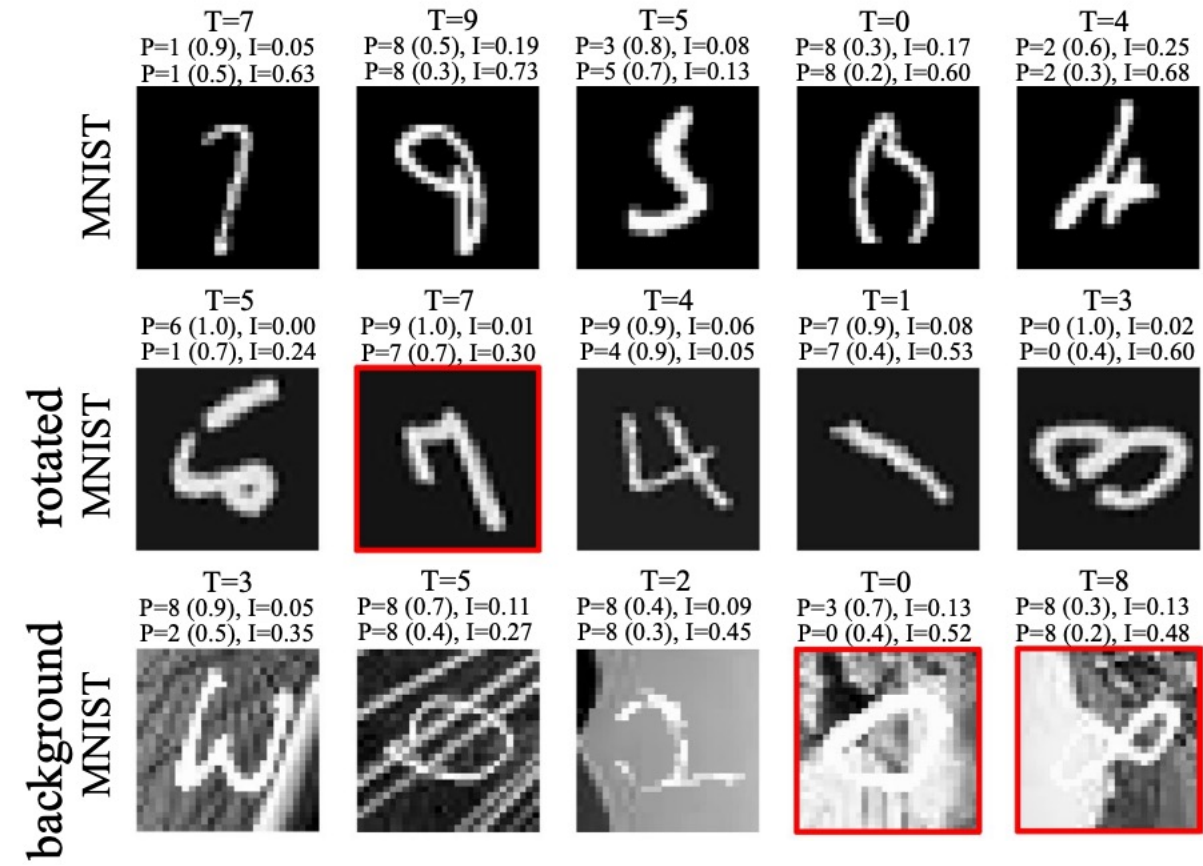# Bayesian meta-learning

HoUston Learning Algorithms (HULA) Lab
Presented by Pengyu (Ben) Yuan

**UNIVERSITY** of **HOUSTON** | **ENGINEERING**

# Uncertainty exists everywhere



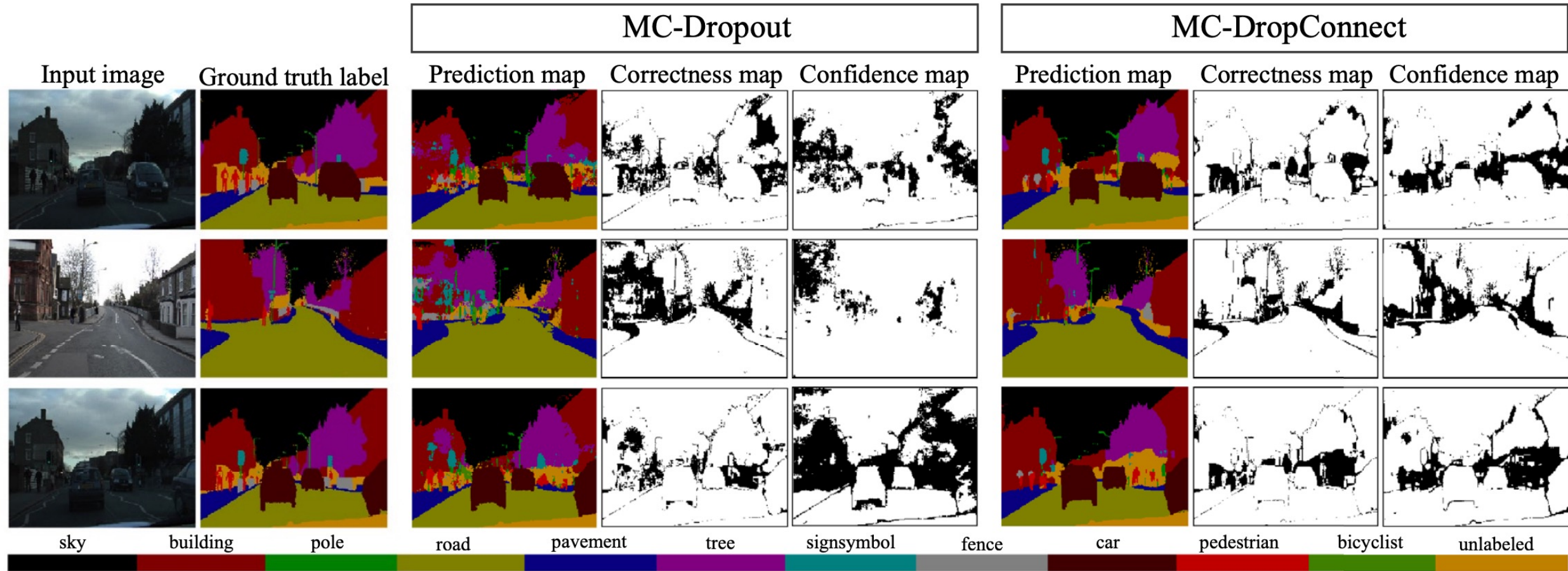**Regression problem: High uncertainty when outside of training distribution**



**Classification problem: Correcting wrong prediction when use uncertainty**

# Uncertainty exists everywhere



**Segmentation problem: High correlation between uncertainty and correctness**

https://www.nature.com/articles/s41598-0__-348_4-x

# Deterministic meta-learning

- Learner's parameter (deterministic)

$$p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right) \approx \delta(\phi_i^*)$$

where

$$\phi_i^* = \arg\max_{\phi_i} \log p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right)$$
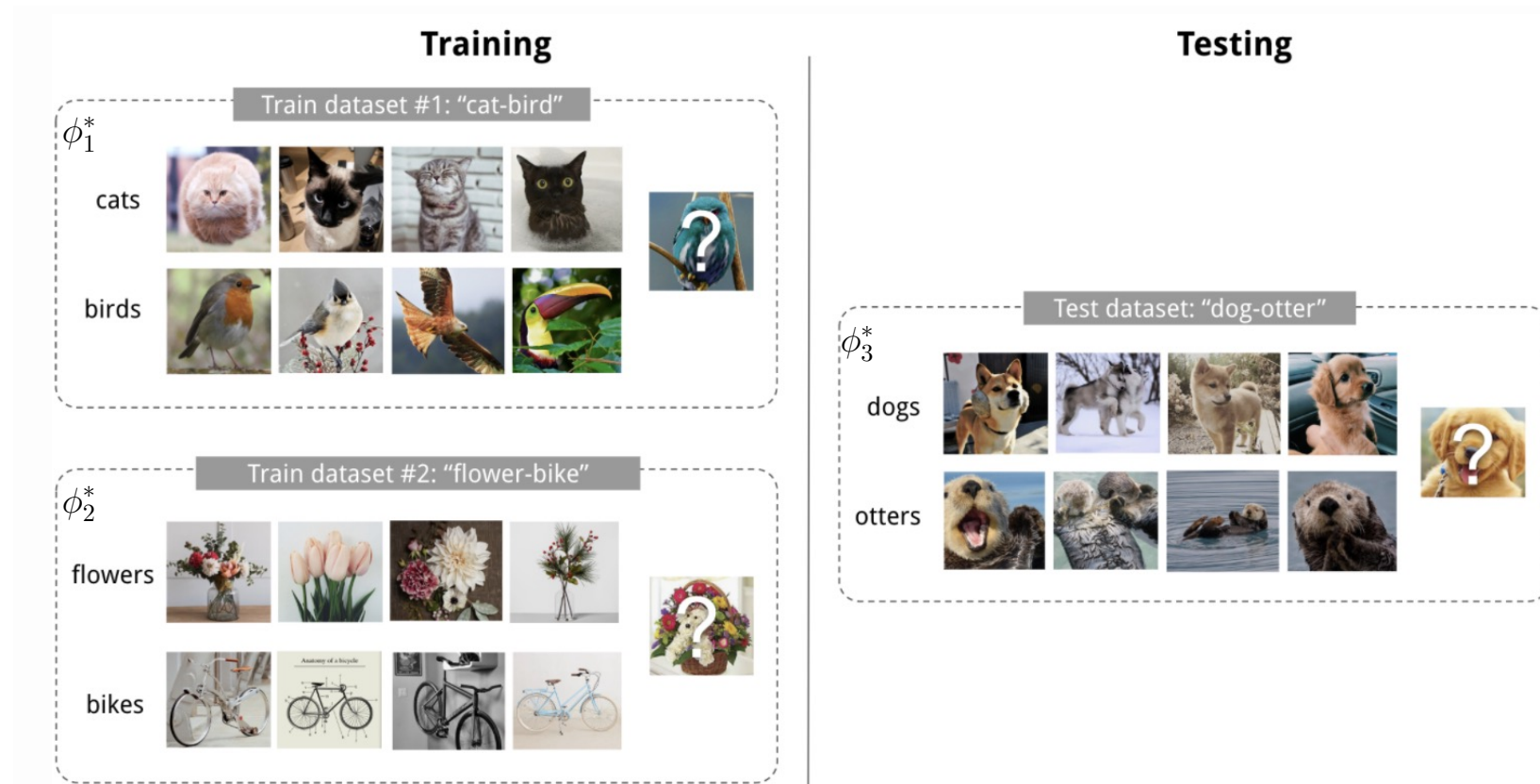
- Problems:
  - $\phi_i^*$ is a point estimate (MAP)
  - There is no uncertainty in the prediction:

    $$y^{\text{te}} = g_{\phi^*}(x^{\text{te}})$$

    where $g$ is the learner's network
  - Few shot learning is ambiguous, easily **overfitting**

Can we get a distribution of $p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right)$ ?

So

$$p(y^{\text{te}}|x^{\text{te}}, \mathcal{D}^{\text{tr}}, \theta^*) = \int_\Phi p(y^{\text{te}}|x^{\text{te}}, \phi)p(\phi|\mathcal{D}^{\text{tr}}, \theta^*)d\phi$$

**Training**

Train dataset #1: "cat-bird"

$\phi_1^*$

cats

birds

Train dataset #2: "flower-bike"

$\phi_2^*$

flowers

bikes

**Testing**

Test dataset: "dog-otter"

$\phi_3^*$

dogs

otters

$\phi_i^*$ **is a set of classifier parameters for** $\mathcal{D}_i$

https://lilianweng.github.io/lil-log/2018/11/30/meta-learning.html

# Deterministic meta-learning

- Meta parameter (deterministic)

$$p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right) \approx \delta(\theta^*)$$

where

$$\theta^* = \arg\max_{\theta} \log p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right)$$
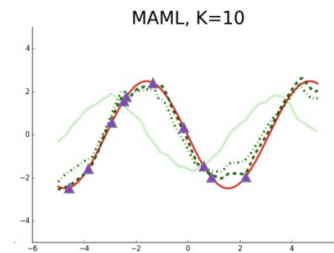
- Problems:
  - $\theta^*$ is also a point estimate (MAP)
  - When the number of tasks is small, there is high uncertainty in the meta parameters. Leads to **meta-overfitting**
  - Learner's parameters are affected by meta parameters:

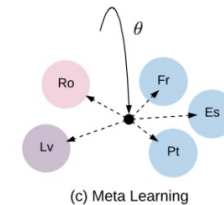$$p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta\right)$$

Thus, it can also affect the distribution of the prediction.

What information might $\boldsymbol{\theta}$ contain…



...in a toy sinusoid problem?

$\boldsymbol{\theta}$ corresponds to family of sinusoid functions (everything but phase and amplitude)

...in multi-language machine translation?

$\boldsymbol{\theta}$ corresponds to the family of all language pairs

$\theta$ **is the shared latent information from** $\mathcal{D}_{\text{meta-train}}$

# Why Bayesian in meta-learning?

Bayesian method can:

- give us a distribution over prediction

- prevent overfitting problem

- update model gradually by using online learning

In meta-learning scenario, it can:

- Learn **safety-critical** few-shot model **(especially in medical imaging)**

- Learn to **actively annotate** new samples **(active learning)**

- Learn to **explore** in meta reinforcement learning

# Bayesian tools

How ?

- Learn distribution over learner's parameters: $\quad \delta(\phi_i^*) \longrightarrow p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right)$

- Learn distribution over meta parameters: $\quad \delta(\theta^*) \longrightarrow p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right)$
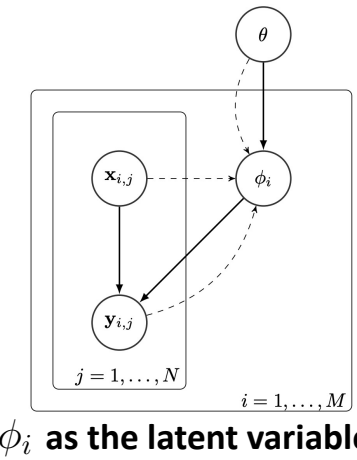
- Can be either one or both

Bayesian toolboxes:

- Latent variable models + variational inference

  approximate likelihood of latent variable model with variational lower bound

- Bayesian ensembles

  particle-based representation: train separate models on bootstraps of the data

- Bayesian neural networks

  explicit distribution over the space of network parameters

- ...

https://openreview.net/pdf?id=rkgpy3C5tX
http://cs330.stanford.edu/fall2020/slides/cs330_bayesian_meta_learning_2020.pdf
https://www.researchgate.net/publication/328757994_A_Batched_Scalable_Multi-Objective_Bayesian_Optimization_Algorithm



$\phi_i$ **as the latent variable**



**Single donkey**          **An ensemble of donkey**



**Bayesian neural network**

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)

- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks

- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
  - Active-learning evaluation

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)

- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks

- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
  - Active-learning evaluation

# The evidence lower bound (ELBO)

- What is ELBO?

It is an optimization function used in variational inference.

$$q(\phi) = \arg\max_q ELBO \qquad \text{where} \qquad ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$

$\phi$ – parameters

$\mathcal{D}$ – observations

$q(\phi)$ – variational distribution

# The evidence lower bound (ELBO)

- **What is ELBO?**

  It is an optimization function used in variational inference.

  $$q(\phi) = \arg \max_q ELBO \qquad \text{where} \qquad ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$

  | | |
  |---|---|
  | $\phi$ | $-$ parameters |
  | $\mathcal{D}$ | $-$ observations |
  | $q(\phi)$ | $-$ variational distribution |

- **What is variational inference (VI)?**

  Approximating a posterior distribution with some easy to manipulate distribution like the Gaussian



$$q(\phi) \to p(\phi|\mathcal{D})$$

# The evidence lower bound (ELBO)

- ## What is ELBO?

  It is an optimization function used in variational inference.

  $$q(\phi) = \arg \max_{q} ELBO \qquad \text{where} \qquad ELBO = \int_{\Phi} q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$
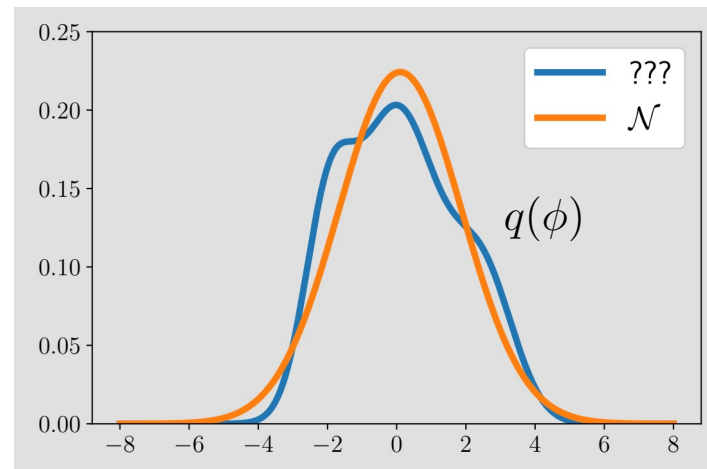
  | | |
  |---|---|
  | $\phi$ | – parameters |
  | $\mathcal{D}$ | – observations |
  | $q(\phi)$ | – variational distribution |

- ## What is variational inference (VI)?

  Approximating a posterior distribution with some easy to manipulate distribution like the Gaussian

  $$q(\phi) \overset{\text{Approx.}}{\rightarrow} p(\phi|\mathcal{D})$$

- ## Why do we need VI?

  It is intractable to calculate the true posterior distribution

  $$p(\phi \mid \mathcal{D}) = \frac{p(\mathcal{D}, \phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{\int_{\Phi} p(\mathcal{D} \mid \phi)p(\phi)d\phi}$$

  because it is impossible to consider all configurations $\phi$ of the neural network.

  | | |
  |---|---|
  | $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
  | $p(\mathcal{D} \mid \phi)$ | – likelihood |
  | $p(\phi)$ | – prior distribution |

# The evidence lower bound (ELBO)

- What is ELBO?

  It is an optimization function used in variational inference.

  $$q(\phi) = \arg\max_q ELBO \qquad \text{where} \qquad ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$
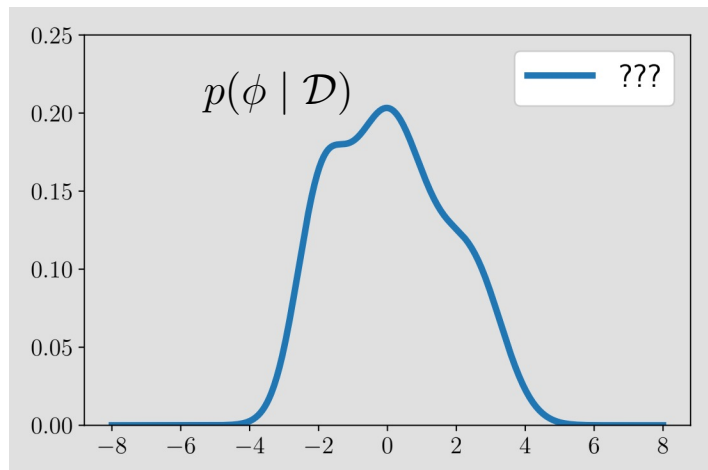
  | $\phi$ | $-$ parameters |
  |---|---|
  | $\mathcal{D}$ | $-$ observations |
  | $q(\phi)$ | $-$ variational distribution |

- What is variational inference (VI)?

  Approximating a posterior distribution with some easy to manipulate distribution like the Gaussian

  $$q(\phi) \overset{\text{Approx.}}{\to} p(\phi|\mathcal{D})$$

- Why do we need VI?

  It is intractable to calculate the true posterior distribution     $p(\phi \mid \mathcal{D}) = \frac{p(\mathcal{D}, \phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{\int_\Phi p(\mathcal{D} \mid \phi)p(\phi)d\phi}$

- Why do we need posterior distribution?

  To get distribution (uncertainty) over predictions:

  $$p(y^{\text{te}}|x^{\text{te}}, \mathcal{D}^{\text{tr}}) = \int_\Phi p(y^{\text{te}}|x^{\text{te}}, \phi)p(\phi|\mathcal{D}^{\text{tr}})d\phi$$

  | $p(\phi \mid \mathcal{D})$ | $-$ true posterior distribution |
  |---|---|
  | $p(\mathcal{D} \mid \phi)$ | $-$ likelihood |
  | $p(\phi)$ | $-$ prior distribution |

# The evidence lower bound (ELBO)

- Why optimizing the ELBO can help to approximate the true posterior distribution?

Look at the Bayesian rule:

Likelihood                                    Prior

Posterior

$$p(\phi \mid \mathcal{D}) = \frac{p(\mathcal{D}, \phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{p(\mathcal{D})}$$

Evidence

$q(\phi)$ – variational distribution

$p(\phi \mid \mathcal{D})$ – true posterior distribution

$p(\mathcal{D} \mid \phi)$ – likelihood

$p(\phi)$ – prior distribution

$p(\mathcal{D})$ – evidence

# The evidence lower bound (ELBO)

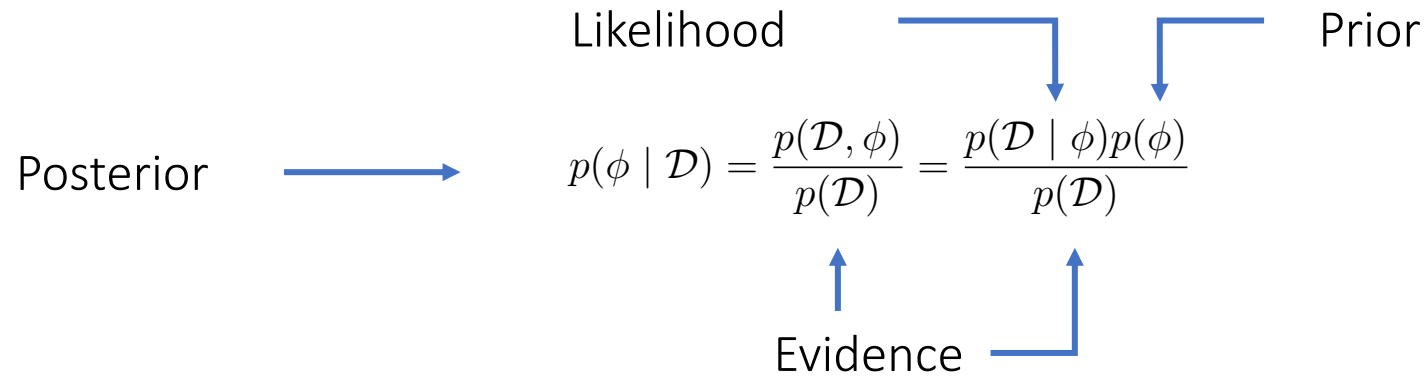- Why optimizing the ELBO can help to approximate the true posterior distribution?

Look at the Bayesian rule:

Likelihood $\qquad\qquad$ Prior

Posterior $\longrightarrow$

$$p(\phi \mid \mathcal{D}) = \frac{p(\mathcal{D}, \phi)}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \phi)p(\phi)}{p(\mathcal{D})}$$

Evidence

Log evidence:

$$\ln p(\mathcal{D}) = ELBO + KL(q(\phi)||p(\phi|\mathcal{D})) \geqslant ELBO \qquad\qquad KL(\cdot||\cdot) \geqslant 0$$

Difference between the variational distribution
and true posterior distribution

$q(\phi)$ – variational distribution

$p(\phi \mid \mathcal{D})$ – true posterior distribution

$p(\mathcal{D} \mid \phi)$ – likelihood

Evidence is fixed by data, thus

$$q(\phi) = \arg\min_q KL(q(\phi)||p(\phi|\mathcal{D})) = \arg\max_q ELBO$$

$p(\phi)$ – prior distribution

$p(\mathcal{D})$ – evidence

# The evidence lower bound (ELBO)

- Let's get a closer look at ELBO

$$\max_q ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$

| | |
|---|---|
| $q(\phi)$ | – variational distribution |
| $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
| $p(\mathcal{D} \mid \phi)$ | – likelihood |
| $p(\phi)$ | – prior distribution |
| $p(\mathcal{D})$ | – evidence |

# The evidence lower bound (ELBO)

- Let's get a closer look at ELBO

$$\max_q ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$

$$= \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}|\phi)p(\phi)}{q(\phi)} d\phi$$

| | |
|---|---|
| $q(\phi)$ | – variational distribution |
| $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
| $p(\mathcal{D} \mid \phi)$ | – likelihood |
| $p(\phi)$ | – prior distribution |
| $p(\mathcal{D})$ | – evidence |

# The evidence lower bound (ELBO)

- Let's get a closer look at ELBO

$$
\begin{aligned}
\max_{q} ELBO &= \int_{\Phi} q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi \\
&= \int_{\Phi} q(\phi) \ln \frac{p(\mathcal{D}|\phi)p(\phi)}{q(\phi)} d\phi \\
&= \int_{\Phi} q(\phi) \ln p(\mathcal{D} \mid \phi) - \int_{\Phi} q(\phi) \ln \frac{q(\phi)}{p(\phi)} d\phi
\end{aligned}
$$

| | |
|---|---|
| $q(\phi)$ | – variational distribution |
| $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
| $p(\mathcal{D} \mid \phi)$ | – likelihood |
| $p(\phi)$ | – prior distribution |
| $p(\mathcal{D})$ | – evidence |

# The evidence lower bound (ELBO)

- Let's get a closer look at ELBO

$$
\begin{aligned}
\max_q ELBO &= \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi \\
&= \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}|\phi)p(\phi)}{q(\phi)} d\phi \\
&= \int_\Phi q(\phi) \ln p(\mathcal{D} \mid \phi) - \int_\Phi q(\phi) \ln \frac{q(\phi)}{p(\phi)} d\phi \\
&= \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))
\end{aligned}
$$

| | |
|---|---|
| $q(\phi)$ | – variational distribution |
| $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
| $p(\mathcal{D} \mid \phi)$ | – likelihood |
| $p(\phi)$ | – prior distribution |
| $p(\mathcal{D})$ | – evidence |

# The evidence lower bound (ELBO)

- Let's get a closer look at ELBO

$$\max_q ELBO = \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}, \phi)}{q(\phi)} d\phi$$

$$= \int_\Phi q(\phi) \ln \frac{p(\mathcal{D}|\phi)p(\phi)}{q(\phi)} d\phi$$

$$= \int_\Phi q(\phi) \ln p(\mathcal{D} \mid \phi) - \int_\Phi q(\phi) \ln \frac{q(\phi)}{p(\phi)} d\phi$$

$$= \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

Samples from $q(\phi)$ to perform original tasks      Regularization term

| | |
|---|---|
| $q(\phi)$ | – variational distribution |
| $p(\phi \mid \mathcal{D})$ | – true posterior distribution |
| $p(\mathcal{D} \mid \phi)$ | – likelihood |
| $p(\phi)$ | – prior distribution |
| $p(\mathcal{D})$ | – evidence |

# The evidence lower bound (ELBO)

- More about ELBO

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence

variational
distribution

likelihood

prior
distribution

# The evidence lower bound (ELBO)

- More about ELBO

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence

variational
distribution

likelihood

prior
distribution

variational distribution can be any form:

$$q(\phi) \Rightarrow q(\phi|\mathcal{D}), \ q(\phi|\theta)$$

For example $\quad q(\phi|\theta) \overset{\text{Approx.}}{\rightarrow} p(\phi|\mathcal{D}) \quad : \quad \ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi|\theta)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi|\theta)\|p(\phi))$

# The evidence lower bound (ELBO)

- More about ELBO

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence      variational distribution      likelihood      prior distribution

variational distribution can be any form:

$$q(\phi) \Rightarrow q(\phi|\mathcal{D}), \ q(\phi|\theta)$$

For example $\quad q(\phi|\theta) \overset{\text{Approx.}}{\to} p(\phi|\mathcal{D}) \quad$ : $\qquad \ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi|\theta)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi|\theta)\|p(\phi))$

posterior can be conditioned on other variables:

$$p(\phi|\mathcal{D}) \Rightarrow p(\phi|\mathcal{D}, \theta)$$

For example $\quad q(\phi) \overset{\text{Approx.}}{\to} p(\phi|\mathcal{D}, \theta) \quad$ : $\qquad \ln p(\mathcal{D}|\theta) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi, \theta)] - KL(q(\phi)\|p(\phi|\theta))$

# The evidence lower bound (ELBO)

- Problem of ELBO:

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence     variational distribution     likelihood     prior distribution

1. Cannot optimizing $q(\phi)$ directly, reparameterization trick is required.
- Variational distribution has variational parameters $\lambda$ , i.e. $q(\phi) = q_\lambda(\phi)$     $\lambda = \{\mu, \sigma^2\}$
- It is in general difficult to calculate the derivative $\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)}$

# The evidence lower bound (ELBO)

- Problem of ELBO:

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi) \| p(\phi))$$

evidence

variational
distribution

likelihood

prior
distribution

1. Cannot optimizing $q(\phi)$ directly, reparameterization trick is required.
- Variational distribution has variational parameters $\lambda$ , i.e. $q(\phi) = q_\lambda(\phi)$ $\lambda = \{\mu, \sigma^2\}$
- It is in general difficult to calculate the derivative $\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)}$

To see that, set $f_\lambda(\phi, \mathcal{D}) = \ln p(\mathcal{D} \mid \phi)$ , then

$$\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)} [f_\lambda(\phi, \mathcal{D})] = \nabla_\lambda \left[ \int_\Phi q_\lambda(\phi) f_\lambda(\phi, \mathcal{D}) d\phi \right]$$

$$= \int_\Phi \nabla_\lambda [q_\lambda(\phi) f_\lambda(\phi, \mathcal{D})] \, d\phi$$

$$= \int_\Phi f_\lambda(\phi, \mathcal{D}) \nabla_\lambda q_\lambda(\phi) d\phi + \int_\Phi q_\lambda(\phi) \nabla_\lambda f_\lambda(\phi, \mathcal{D}) d\phi$$

$$= \underbrace{\int_\Phi f_\lambda(\phi, \mathcal{D}) \nabla_\lambda q_\lambda(\phi) d\phi}_{\text{What about this?}} + \mathbb{E}_{q_\lambda(\phi)} [\nabla_\lambda f_\lambda(\phi, \mathcal{D})]$$

https://gregorygundersen.com/blog/2018/04/29/reparameterization/

# The evidence lower bound (ELBO)

- Problem of ELBO:

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence     variational distribution     likelihood     prior distribution

1. Cannot optimizing $q(\phi)$ directly, reparameterization trick is required.
- Variational distribution has variational parameters $\lambda$ , i.e. $q(\phi) = q_\lambda(\phi)$     $\lambda = \{\mu, \sigma^2\}$
- It is in general difficult to calculate the derivative $\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)}$

To see that, set $f_\lambda(\phi, \mathcal{D}) = \ln p(\mathcal{D} \mid \phi)$ , then

reparameterization

$$\phi \sim N(\mu, \sigma^2)$$

$$\phi = \mu + \sigma\epsilon, \text{ where } \epsilon \sim N(0, \mathrm{I})$$

$$f_\lambda(\phi, \mathcal{D}) = f(g_\lambda(\epsilon, \mathcal{D}))$$

$$\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)}[f_\lambda(\phi, \mathcal{D})] = \underbrace{\int_\Phi f_\lambda(\phi, \mathcal{D}) \nabla_\lambda q_\lambda(\phi) d\phi}_{\text{What about this?}} + \mathbb{E}_{q_\lambda(\phi)}[\nabla_\lambda f_\lambda(\phi, \mathcal{D})]$$

$$\nabla_\lambda \mathbb{E}_{q_\lambda(\phi)}[f_\lambda(\phi, \mathcal{D})] = \nabla_\lambda \mathbb{E}_{p(\epsilon)}[f(g_\lambda(\epsilon, \mathcal{D}))]$$
$$= \mathbb{E}_{p(\epsilon)}[\nabla_\lambda f(g_\lambda(\epsilon, \mathcal{D}))]$$

# The evidence lower bound (ELBO)

- Problem of ELBO:

$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

evidence

variational distribution

likelihood

prior distribution

2. **Can only model Gaussian variational distribution** $q(\phi)$

- Variational distribution needs to be simple
- Reparameterization trick gives us Gaussian distribution
- KL divergence has analytic solution when both distributions are Gaussian

# The evidence lower bound (ELBO)

Basic ideas: maximize ELBO to use $q(\phi)$ to approximate $p(\phi \mid \mathcal{D})$

- During training:
    1. Sample model parameters from $q(\phi)$
    2. Maximize the likelihood of the observation $p(\mathcal{D} \mid \phi)$ while minimize the gap between the $q(\phi)$ and $p(\phi)$

$$\max \ \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

- During test:
    1. Sample model parameters from $q(\phi)$
    2. Use it as the true posterior distribution for prediction

$$p(y^{\text{te}}|x^{\text{te}}, \mathcal{D}^{\text{tr}}) = \int_{\Phi} p(y^{\text{te}}|x^{\text{te}}, \phi)p(\phi|\mathcal{D}^{\text{tr}})d\phi$$
$$\approx \int_{\Phi} p(y^{\text{te}}|x^{\text{te}}, \phi)q(\phi)d\phi$$

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)
- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks
- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
  - Active-learning evaluation

# Amortized variational inference

For dataset $\mathcal{D}_i$ , the posterior distribution we need is: $p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta\right)$

The variational distribution we used: $q\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta\right) \sim N(\mu_i, \sigma_i^2)$

## Amortized Variational Inference

1. Introduce a parameterized model or function that outputs the variational parameters of the approximate posteriors.

$$\lambda = \left\{\mu_i, \sigma_i^2\right\} = f_\theta(\mathcal{D}_i^{\mathrm{tr}})$$

$$\lambda = \left\{\mu_i, \sigma_i^2\right\} = f(\mathcal{D}_i^{\mathrm{tr}}, \theta)$$



2. Variational parameter $\lambda$ is determined by meta parameter $\theta$, thus optimizing variational parameter is the same as optimizing the meta parameter. This optimization is done by doing gradient descent on the **loss function for the variational inference**.

# Amortized variational inference

- What is the loss function for the variational inference in meta-learning?

Variational inference:
$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

Because of the meta-learning, we have additional meta-parameter $\theta$



By replacing:

$$q(\phi) \Rightarrow q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta) \qquad p(\phi|\mathcal{D}) \Rightarrow p(\phi|\mathcal{D}, \theta)$$

We have:

$$\max_{\theta} ELBO = \max_{\theta} \mathbb{E}_{q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)}[\ln p(\mathcal{D}|\phi, \theta)] - KL(q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)\|p(\phi|\theta))$$

# Amortized variational inference

- What is the loss function for the variational inference in meta-learning?

Variational inference:
$$\ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$$

Because of the meta-learning, we have additional meta-parameter $\theta$



By replacing: $\qquad q(\phi) \Rightarrow q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta) \qquad\qquad p(\phi|\mathcal{D}) \Rightarrow p(\phi|\mathcal{D}, \theta)$

We have: $\qquad \max_\theta ELBO = \max_\theta \mathbb{E}_{q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)}[\ln p(\mathcal{D}|\phi, \theta)] - KL(q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)\|p(\phi|\theta))$

Approximating the posterior of test data: $\qquad p(\phi|\mathcal{D}, \theta) \Rightarrow p(\phi|\mathcal{D}^{\mathrm{te}}, \theta)$

We have: $\qquad \max_\theta ELBO = \max_\theta \mathbb{E}_{q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)}[\ln p(\underline{\mathcal{D}^{\mathrm{te}}}|\phi, \theta)] - KL(q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)\|p(\phi|\theta))$
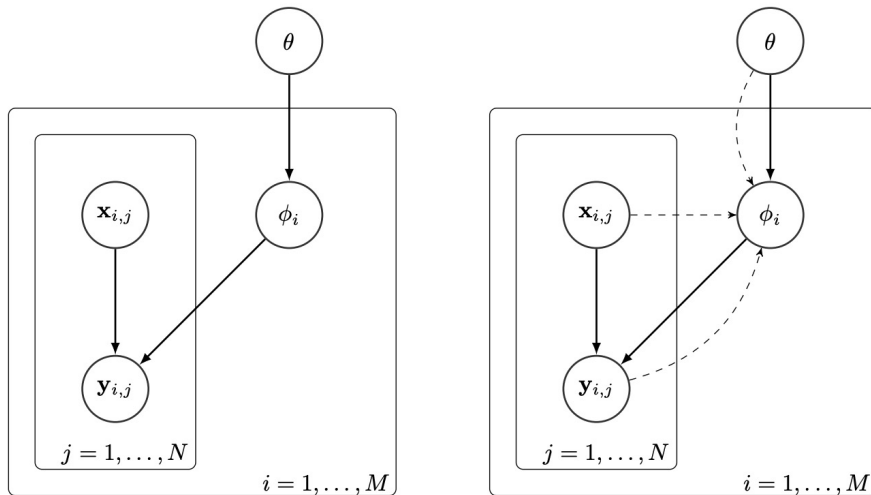
# Amortized variational inference

- What is the loss function for the variational inference in meta-learning?

Variational inference: $\quad \ln p(\mathcal{D}) \geqslant ELBO = \mathbb{E}_{q(\phi)}[\ln p(\mathcal{D} \mid \phi)] - KL(q(\phi)\|p(\phi))$

Because of the meta-learning, we have additional meta-parameter $\theta$
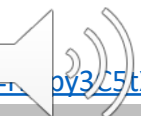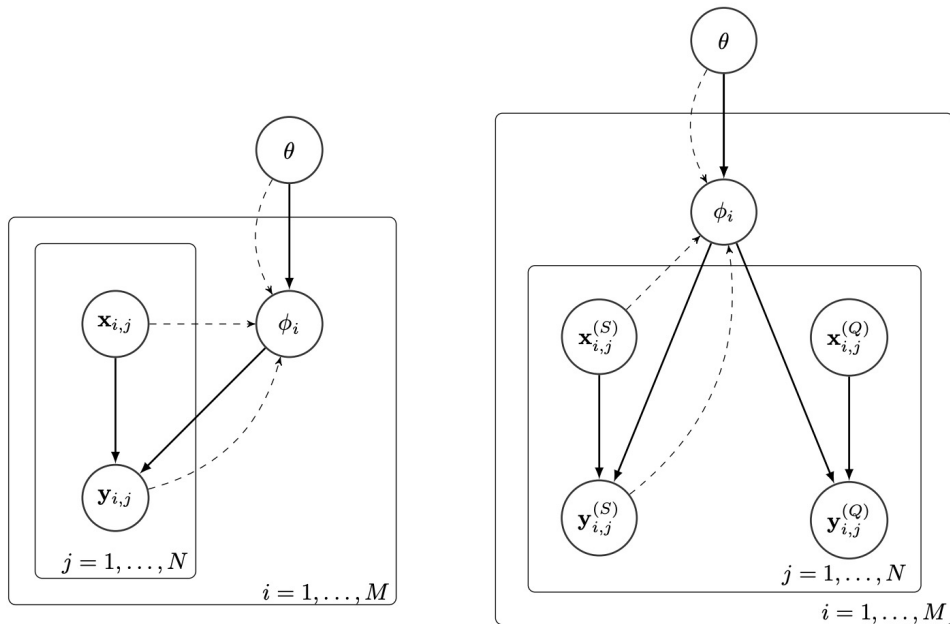
By replacing: $\qquad q(\phi) \Rightarrow q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta) \qquad\qquad p(\phi|\mathcal{D}) \Rightarrow p(\phi|\mathcal{D}, \theta)$

We have: $\qquad \max_{\theta} ELBO = \max_{\theta} \mathbb{E}_{q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)}[\ln p(\mathcal{D}|\phi, \theta)] - KL(q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)\|p(\phi|\theta))$

Approximating the posterior of test data: $\qquad p(\phi|\mathcal{D}, \theta) \Rightarrow p(\phi|\mathcal{D}^{\mathrm{te}}, \theta)$

We have: $\qquad \max_{\theta} ELBO = \max_{\theta} \mathbb{E}_{q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)}[\ln p(\underline{\mathcal{D}^{\mathrm{te}}}|\phi, \theta)] - KL(q(\phi|\mathcal{D}^{\mathrm{tr}}, \theta)\|p(\phi|\theta))$

For all tasks, the final objective is:

$$\max_{\theta} \; \underline{\mathbb{E}_{p(\mathcal{D}_i)}} \left[ \underline{\mathbb{E}_{q(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta)}[\ln p(\mathcal{D}_i^{\mathrm{te}}|\phi_i, \theta)] - KL(q(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta)\|p(\phi_i|\theta))} \right]$$

$ELBO$ for $\mathcal{D}_i$

# Amortized variational inference

**Different way to model** $q\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta\right) \sim N(\mu_i, \sigma_i^2)$

$$\max_\theta \ \mathbb{E}_{p(\mathcal{D}_i)} \left[ \mathbb{E}_{q(\phi_i|\mathcal{D}_i^{\mathrm{tr}},\theta)}[\ln p(\mathcal{D}_i^{\mathrm{te}}|\phi_i,\theta)] - KL(q(\phi_i|\mathcal{D}_i^{\mathrm{tr}},\theta)\|p(\phi_i|\theta)) \right]$$

## Two parametric meta-learning approaches:

|  | Deterministic version | Bayesian version |
|---|---|---|

- Black-box based.

  **Key idea:** Train a neural network to represent
  $$q\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta\right)$$

  Deterministic version: $\phi_i = f_\theta(\mathcal{D}_i^{\mathrm{tr}})$

  Bayesian version: $\lambda = \{\mu_i, \sigma_i^2\} = f(\mathcal{D}_i^{\mathrm{tr}}, \theta)$

- Optimization based (normally using gradient descent).

  **Key idea:** Acquire $\phi_i$ through optimization on $\mathcal{D}_i^{\mathrm{tr}}$,

  meta parameter $\theta$ serves as a prior

  Deterministic version: $\phi_i = \arg\max_{\phi_i} \log p\left(\mathcal{D}_i^{\mathrm{tr}} \mid \phi_i\right) + \log p\left(\phi_i \mid \theta\right)$

  Bayesian version: $\lambda_i = \arg\max_{\lambda_i} \log p\left(\mathcal{D}_i^{\mathrm{tr}} \mid \phi_i\right) + \log p\left(\phi_i \mid \theta\right)$

# Amortized variational inference

- Black-box based approach (VERSA)

$$\lambda = \{\mu_i, \sigma_i^2\} = f(\mathcal{D}_i^{\mathrm{tr}}, \theta)$$

- Meta learner has two components:
  - Feature extraction network $\theta$
  - Amortization network $\phi$

- Meta parameter $\{\theta, \phi\}$

- Learner's parameter $\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$



Training data $D^{(t)} = \{(x_n^{(t)}, y_n^{(t)})\}_{n=1}^{N_t}$, and test data $\{(\tilde{x}_m^{(t)}, \tilde{y}_m^{(t)})\}_{m=1}^{M_t}$

task specific parameters $\{\psi^{(t)}\}_{t=1}^{T}$ $\qquad$ $\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$

# Amortized variational inference

- Black-box based approach (VERSA)

$$\lambda = \{\mu_i, \sigma_i^2\} = f(\mathcal{D}_i^{\text{tr}}, \theta)$$

- Meta learner has two components:
  - Feature extraction network $\theta$
  - Amortization network $\phi$

- Meta parameter $\{\theta, \phi\}$

- Learner's parameter $\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$

- Approximate the posterior with

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$



$$\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$$

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

Collect $L$ samples of $\psi$

# Amortized variational inference
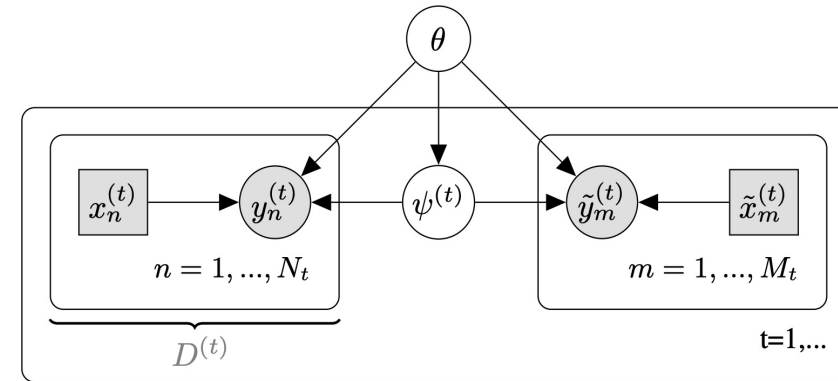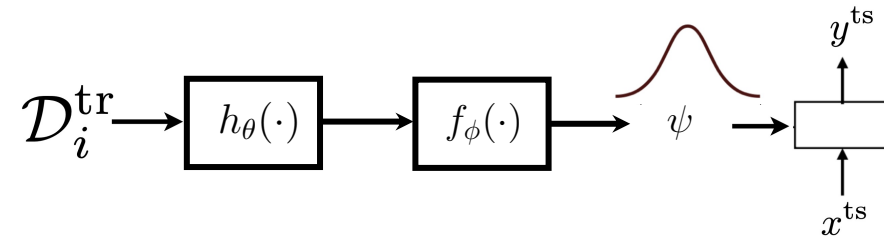
- Black-box based approach (VERSA)

$$\lambda = \{\mu_i, \sigma_i^2\} = f(\mathcal{D}_i^{\text{tr}}, \theta)$$

- Meta learner has two components:
  - Feature extraction network $\theta$
  - Amortization network $\phi$

- Meta parameter $\{\theta, \phi\}$

- Learner's parameter $\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$

- Approximate the posterior with

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

- Objective function

$$\hat{\mathcal{L}}(\theta, \phi) = \frac{1}{MT} \sum_{M,T} \log \frac{1}{L} \sum_{l=1}^{L} p\left(\tilde{y}_m^{(t)} | \tilde{x}_m^{(t)}, \psi_l^{(t)}, \theta\right), \quad \text{with } \psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

$$\max_\theta \mathbb{E}_{p(\mathcal{D}_i)} \left[ \mathbb{E}_{q(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)} [\ln p(\mathcal{D}_i^{\text{te}} | \phi_i, \theta)] - KL(q(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta) \| p(\phi_i | \theta)) \right]$$



Feature extraction — Linear Classifier — Softmax output

$$\tilde{x} \rightarrow h_\theta(\tilde{x}) \rightarrow \left(\begin{array}{ccc} | & & | \\ w_t^{(1)} & \cdots & w_t^{(C)} \\ | & & | \end{array}\right) \rightarrow p(\tilde{y} | \tilde{x}, \theta, \psi_t)$$

$$\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$$

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

Collect $L$ samples of $\psi$

Amortization Network

$$h_\theta\left(x_1^{(1)}\right) \cdots h_\theta\left(x_{k_1}^{(1)}\right) \quad h_\theta\left(x_1^{(C)}\right) \cdots h_\theta\left(x_{k_C}^{(C)}\right)$$

$k_1$ train examples from class 1    $k_C$ train examples from class $C$

$$\mathcal{D}_i^{\text{tr}} \rightarrow h_\theta(\cdot) \rightarrow f_\phi(\cdot) \rightarrow \psi \rightarrow \boxed{} \rightarrow y^{\text{ts}}, x^{\text{ts}}$$

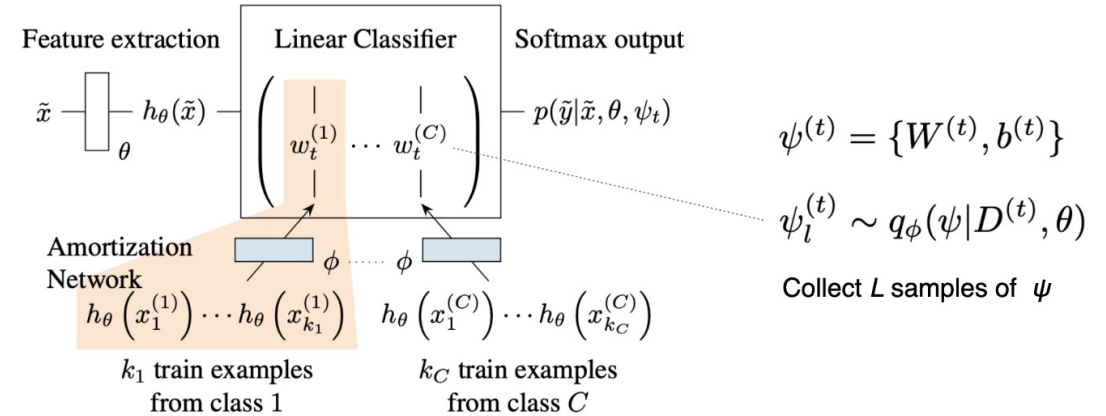| $\delta(\phi_i^*)$ | $p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right)$ |
| $\delta(\theta^*)$ | $p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right)$ |

# Amortized variational inference

- Black-box based approach (VERSA)

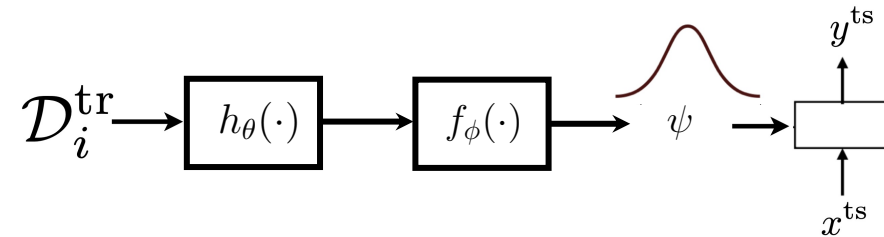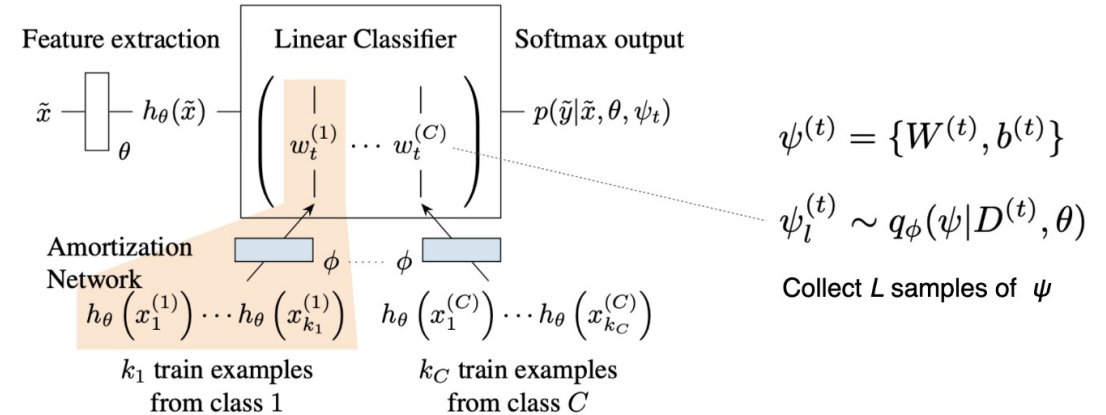$$\lambda = \{\mu_i, \sigma_i^2\} = f(\mathcal{D}_i^{\mathrm{tr}}, \theta)$$

- Meta learner has two components:
  - Feature extraction network $\theta$
  - Amortization network $\phi$

- Meta parameter $\{\theta, \phi\}$

- Learner's parameter $\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$

- Approximate the posterior with

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

- Objective function

$$\hat{\mathcal{L}}(\theta, \phi) = \frac{1}{MT} \sum_{M,T} \log \frac{1}{L} \sum_{l=1}^{L} p\left(\tilde{y}_m^{(t)} | \tilde{x}_m^{(t)}, \psi_l^{(t)}, \theta\right), \quad \text{with } \psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

$$\max_\theta \mathbb{E}_{p(\mathcal{D}_i)} \left[ \mathbb{E}_{q(\phi_i | \mathcal{D}_i^{\mathrm{tr}}, \theta)}[\ln p(\mathcal{D}_i^{\mathrm{te}} | \phi_i, \theta)] - KL(q(\phi_i | \mathcal{D}_i^{\mathrm{tr}}, \theta) \| p(\phi_i | \theta)) \right]$$



$$\psi^{(t)} = \{W^{(t)}, b^{(t)}\}$$

$$\psi_l^{(t)} \sim q_\phi(\psi | D^{(t)}, \theta)$$

Collect $L$ samples of $\psi$



| $\delta(\phi_i^*)$ | ✗ | $p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right)$ | ✓ |
| $\delta(\theta^*)$ | ✓ | $p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right)$ | ✗ |

# Amortized variational inference

- Optimization based approach (Amortized Bayesian Meta-Learning)

- Recall objective function:

$$\max_{\theta} \; \mathbb{E}_{p(\mathcal{D}_i)} \left[ \mathbb{E}_{q(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta)} [\ln p(\mathcal{D}_i^{\text{te}} | \phi_i, \theta)] - KL(q(\phi_i | \mathcal{D}_i^{\text{tr}}, \theta) \| p(\phi_i | \theta)) \right]$$

- $q\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta\right)$ is achieved by SGD on the mean and variance using $\mathcal{D}_i^{\text{tr}}$

$$\lambda_i = \arg \max_{\lambda_i} \log p\left(\mathcal{D}_i^{\text{tr}} \mid \phi_i\right) + \log p\left(\phi_i \mid \theta\right)$$

$$q_\theta \left(\phi_i \middle| \mathcal{D}_i^{(S)}\right) = \mathcal{N}\left(\phi_i; \boldsymbol{\mu}_\lambda^{(K)}, \boldsymbol{\sigma^2}_\lambda^{(K)}\right)$$

- SGD:

1. $\lambda^{(0)} = \lambda^{(init)}$
2. **for** $k = 0, \ldots, K - 1$, **set**
   $\lambda^{(k+1)} = \lambda^{(k)} - \alpha \nabla_{\lambda^{(k)}} \mathcal{L}_{\mathcal{D}}(\lambda^{(k)}, \theta)$

- $\lambda^{(init)}$ is $\theta$ like MAML, SGD is the inner loop optimization

# Amortized variational inference

- Optimization based approach (Amortized Bayesian Meta-Learning)

$$\delta(\phi_i^*) \qquad p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right)$$

$$\delta(\theta^*) \qquad p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right)$$

**Algorithm 1** Meta-training

**Input**: Number of update steps $K$, Number of total episodes $M$, Inner learning rate $\alpha$, Outer learning rate $\beta$

1: Initialize $\theta = \{\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta^2\}$
2: $p(\theta) = \mathcal{N}(\boldsymbol{\mu}; \mathbf{0}, \mathbf{I}) \cdot \prod_{l=1}^{D} \mathrm{Gamma}(\tau_l; a_0, b_0)$
3: **for** $i = 1$ to $M$ **do**
4: $\quad \mathcal{D}_i = \left\{\mathcal{D}_i^{(S)}, \mathcal{D}_i^{(Q)}\right\}$
5: $\quad \boldsymbol{\mu}_\lambda^{(0)} \leftarrow \boldsymbol{\mu}_\theta; \boldsymbol{\sigma}_\lambda^{2(0)} \leftarrow \boldsymbol{\sigma}_\theta^2$
6: $\quad$ **for** $k = 0$ to $K - 1$ **do**
7: $\qquad \lambda^{(k)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(k)}, \boldsymbol{\sigma}_\lambda^{(k)}\right\}$
8: $\qquad \boldsymbol{\mu}_\lambda^{(k+1)} \leftarrow \boldsymbol{\mu}_\lambda^{(k)} - \alpha \nabla_{\boldsymbol{\mu}_\lambda^{(k)}} \mathcal{L}_{\mathcal{D}_i^{(S)}}\left(\lambda^{(k)}, \theta\right)$
9: $\qquad \boldsymbol{\sigma}_\lambda^{2(k+1)} \leftarrow \boldsymbol{\sigma}_\lambda^{2(k)} - \alpha \nabla_{\boldsymbol{\sigma}_\lambda^{2(k)}} \mathcal{L}_{\mathcal{D}_i^{(S)}}\left(\lambda^{(k)}, \theta\right)$
10: $\quad$ **end for**
11:
12: $\quad \lambda^{(K)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(K)}, \boldsymbol{\sigma}_\lambda^{2(K)}\right\}$
13: $\quad q(\theta) = \mathbb{1}\{\boldsymbol{\mu} = \boldsymbol{\mu}_\theta\} \cdot \mathbb{1}\{\boldsymbol{\sigma}^2 = \boldsymbol{\sigma}_\theta^2\}$
14: $\quad \boldsymbol{\mu}_\theta \leftarrow \boldsymbol{\mu}_\theta - \beta \nabla_{\boldsymbol{\mu}_\theta}\left[\mathcal{L}_{\mathcal{D}_i}(\lambda^{(K)}, \theta) + \frac{1}{M}\mathrm{KL}(q(\theta)\|p(\theta))\right]$
15: $\quad \boldsymbol{\sigma}_\theta^2 \leftarrow \boldsymbol{\sigma}_\theta^2 - \beta \nabla_{\boldsymbol{\sigma}_\theta^2}\left[\mathcal{L}_{\mathcal{D}_i}(\lambda^{(K)}, \theta) + \frac{1}{M}\mathrm{KL}(q(\theta)\|p(\theta))\right]$
16: **end for**

SGD

**Algorithm 2** Meta-evaluation

**Input**: Number of update steps $K$, Dataset $\mathcal{D} = \{\mathcal{D}^{(S)}, \mathcal{D}^{(Q)}\}$, Parameters $\theta = \{\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta^2\}$, Inner learning rate $\alpha$

1: $\boldsymbol{\mu}_\lambda^{(0)} \leftarrow \boldsymbol{\mu}_\theta; \boldsymbol{\sigma}_\lambda^{2(0)} \leftarrow \boldsymbol{\sigma}_\theta^2$
2: **for** $k = 0$ to $K - 1$ **do**
3: $\quad \lambda^{(k)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(k)}, \boldsymbol{\sigma}_\lambda^{(k)}\right\}$
4: $\quad \boldsymbol{\mu}_\lambda^{(k+1)} \leftarrow \boldsymbol{\mu}_\lambda^{(k)} - \alpha \nabla_{\boldsymbol{\mu}_\lambda^{(k)}} \mathcal{L}_{\mathcal{D}^{(S)}}\left(\lambda^{(k)}, \theta\right)$
5: $\quad \boldsymbol{\sigma}_\lambda^{2(k+1)} \leftarrow \boldsymbol{\sigma}_\lambda^{2(k)} - \alpha \nabla_{\boldsymbol{\sigma}_\lambda^{2(k)}} \mathcal{L}_{\mathcal{D}^{(S)}}\left(\lambda^{(k)}, \theta\right)$
6: **end for**
7:
8: $q_\theta\left(\phi \mid D^{(S)}\right) = \mathcal{N}\left(\phi; \boldsymbol{\mu}_\lambda^{(K)}, \boldsymbol{\sigma}_\lambda^{2(K)}\right)$
9: Evaluate $D^{(Q)}$ using $\mathbb{E}_{q_\theta(\phi \mid D^{(S)})}\left[p(D^{(Q)} \mid \phi)\right]$

https://openreview.net/pdf?id=rkgpy3C5tX
https://jonathan-hui.medium.com/meta-learning-bayesian-meta-learning-weak-supervision-to-09b2eff3

# Amortized variational inference

- Optimization based approach (Amortized Bayesian Meta-Learning)

$$\delta(\phi_i^*) \quad \text{✗} \quad p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right) \quad \text{✓}$$

$$\delta(\theta^*) \quad \text{✗} \quad p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right) \quad \text{✓}$$

**Algorithm 1** Meta-training

**Input**: Number of update steps $K$, Number of total episodes $M$, Inner learning rate $\alpha$, Outer learning rate $\beta$

1: Initialize $\theta = \{\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta^2\}$
2: $p(\theta) = \mathcal{N}(\boldsymbol{\mu}; \mathbf{0}, \mathbf{I}) \cdot \prod_{l=1}^{D} \mathrm{Gamma}(\tau_l; a_0, b_0)$
3: **for** $i = 1$ to $M$ **do**
4: $\quad \mathcal{D}_i = \left\{\mathcal{D}_i^{(S)}, \mathcal{D}_i^{(Q)}\right\}$
5: $\quad \boldsymbol{\mu}_\lambda^{(0)} \leftarrow \boldsymbol{\mu}_\theta; \boldsymbol{\sigma}_\lambda^{2(0)} \leftarrow \boldsymbol{\sigma}_\theta^2$
6: $\quad$ **for** $k = 0$ to $K - 1$ **do**
7: $\qquad \lambda^{(k)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(k)}, \boldsymbol{\sigma}_\lambda^{(k)}\right\}$
8: $\qquad \boldsymbol{\mu}_\lambda^{(k+1)} \leftarrow \boldsymbol{\mu}_\lambda^{(k)} - \alpha \nabla_{\boldsymbol{\mu}_\lambda^{(k)}} \mathcal{L}_{\mathcal{D}_i^{(S)}}\left(\lambda^{(k)}, \theta\right)$
9: $\qquad \boldsymbol{\sigma}_\lambda^{2(k+1)} \leftarrow \boldsymbol{\sigma}_\lambda^{2(k)} - \alpha \nabla_{\boldsymbol{\sigma}_\lambda^{2(k)}} \mathcal{L}_{\mathcal{D}_i^{(S)}}\left(\lambda^{(k)}, \theta\right)$
10: $\quad$ **end for**
11:
12: $\quad \lambda^{(K)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(K)}, \boldsymbol{\sigma}_\lambda^{2(K)}\right\}$
13: $\quad q(\theta) = \mathbb{1}\{\boldsymbol{\mu} = \boldsymbol{\mu}_\theta\} \cdot \mathbb{1}\{\boldsymbol{\sigma}^2 = \boldsymbol{\sigma}_\theta^2\}$
14: $\quad \boldsymbol{\mu}_\theta \leftarrow \boldsymbol{\mu}_\theta - \beta \nabla_{\boldsymbol{\mu}_\theta}\left[\mathcal{L}_{\mathcal{D}_i}(\lambda^{(K)}, \theta) + \frac{1}{M}\mathrm{KL}(q(\theta)\|p(\theta))\right]$
15: $\quad \boldsymbol{\sigma}_\theta^2 \leftarrow \boldsymbol{\sigma}_\theta^2 - \beta \nabla_{\boldsymbol{\sigma}_\theta^2}\left[\mathcal{L}_{\mathcal{D}_i}(\lambda^{(K)}, \theta) + \frac{1}{M}\mathrm{KL}(q(\theta)\|p(\theta))\right]$
16: **end for**

SGD

**Algorithm 2** Meta-evaluation

**Input**: Number of update steps $K$, Dataset $\mathcal{D} = \{\mathcal{D}^{(S)}, \mathcal{D}^{(Q)}\}$, Parameters $\theta = \{\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta^2\}$, Inner learning rate $\alpha$

1: $\boldsymbol{\mu}_\lambda^{(0)} \leftarrow \boldsymbol{\mu}_\theta; \boldsymbol{\sigma}_\lambda^{2(0)} \leftarrow \boldsymbol{\sigma}_\theta^2$
2: **for** $k = 0$ to $K - 1$ **do**
3: $\quad \lambda^{(k)} \leftarrow \left\{\boldsymbol{\mu}_\lambda^{(k)}, \boldsymbol{\sigma}_\lambda^{(k)}\right\}$
4: $\quad \boldsymbol{\mu}_\lambda^{(k+1)} \leftarrow \boldsymbol{\mu}_\lambda^{(k)} - \alpha \nabla_{\boldsymbol{\mu}_\lambda^{(k)}} \mathcal{L}_{\mathcal{D}^{(S)}}\left(\lambda^{(k)}, \theta\right)$
5: $\quad \boldsymbol{\sigma}_\lambda^{2(k+1)} \leftarrow \boldsymbol{\sigma}_\lambda^{2(k)} - \alpha \nabla_{\boldsymbol{\sigma}_\lambda^{2(k)}} \mathcal{L}_{\mathcal{D}^{(S)}}\left(\lambda^{(k)}, \theta\right)$
6: **end for**
7:
8: $q_\theta\left(\phi \mid D^{(S)}\right) = \mathcal{N}\left(\phi; \boldsymbol{\mu}_\lambda^{(K)}, \boldsymbol{\sigma}_\lambda^{2(K)}\right)$
9: Evaluate $D^{(Q)}$ using $\mathbb{E}_{q_\theta(\phi \mid D^{(S)})}\left[p(D^{(Q)} \mid \phi)\right]$

https://openreview.net/pdf?id=rkgpy3C5tX
https://jonathan-hui.medium.com/meta-learning-bayesian-meta-learning-weak-supervision-10910b2ff3
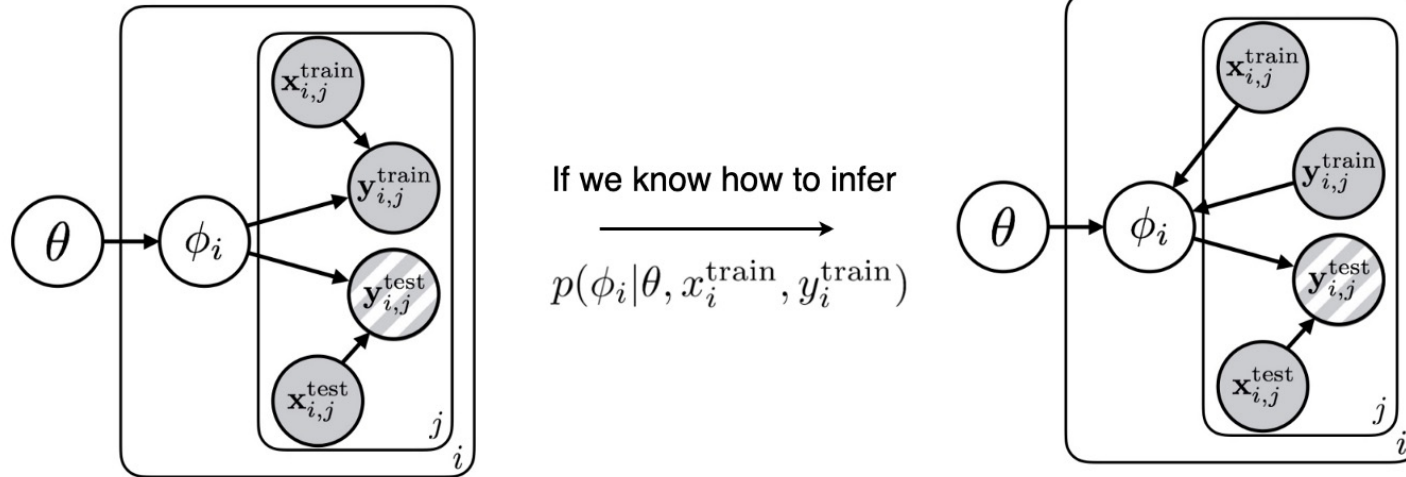
# Amortized variational inference

- Optimization based approach (Probabilistic MAML)

How about approximating posterior distribution directly on $\theta$ ?

$$\delta(\phi_i^*) \quad \checkmark \qquad p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right) \qquad \mathsf{X}$$

$$\delta(\theta^*) \quad \mathsf{X} \qquad p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right) \quad \checkmark$$



Original graphical model

If we know how to infer

$$p(\phi_i \mid \theta, x_i^{\mathrm{train}}, y_i^{\mathrm{train}})$$

For example, in MAML,

$p(\phi_i \mid \mathbf{x}_i^{\mathrm{tr}}, \mathbf{y}_i^{\mathrm{tr}}, \theta) \approx \delta(\phi_i = \phi_i^\star)$ where $\phi_i^\star$ is obtained via gradient descent starting from $\theta$.

# Amortized variational inference

- Optimization based approach (Probabilistic MAML)

$$\theta \sim p(\theta) = \mathcal{N}(\mu_\theta, \Sigma_\theta)$$

**key idea:** $p(\phi_i | \theta, x_i^{\text{train}}, y_i^{\text{train}}) \approx \delta(\hat{\phi}_i)$ $\qquad \hat{\phi}_i \approx \theta + \alpha \nabla_\theta \log p(y_i^{\text{train}} | x_i^{\text{train}}, \theta)$

**What does ancestral sampling look like?**

1. $\theta \sim \mathcal{N}(\mu_\theta, \Sigma_\theta)$
2. $\phi_i \sim p(\phi_i | \theta, x_i^{\text{train}}, y_i^{\text{train}}) \approx \hat{\phi}_i = \theta + \alpha \nabla_\theta \log p(y_i^{\text{train}} | x_i^{\text{train}}, \theta)$

$\mathcal{L}(\phi, \mathcal{D}_{\text{train}})$

$\mu_\theta$

smiling, hat

smiling, young

$\phi$

# Amortized variational inference

- Optimization based approach (Probabilistic MAML)

**Algorithm 1** Meta-training, differences from MAML in red

**Require:** $p(\mathcal{T})$: distribution over tasks
1: initialize $\Theta := \{\boldsymbol{\mu}_\theta, \boldsymbol{\sigma}_\theta^2, \mathbf{v}_q, \boldsymbol{\gamma}_p, \boldsymbol{\gamma}_q\}$
2: **while** not done **do**
3:     Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
4:     **for all** $\mathcal{T}_i$ **do**
5:         $\mathcal{D}^{\text{tr}}, \mathcal{D}^{\text{test}} = \mathcal{T}_i$
6:         Evaluate $\nabla_{\boldsymbol{\mu}_\theta} \mathcal{L}(\boldsymbol{\mu}_\theta, \mathcal{D}^{\text{test}})$
7:         Sample $\theta \sim q = \mathcal{N}(\boldsymbol{\mu}_\theta - \boldsymbol{\gamma}_q \nabla_{\boldsymbol{\mu}_\theta} \mathcal{L}(\boldsymbol{\mu}_\theta, \mathcal{D}^{\text{test}}), \mathbf{v}_q)$
8:         Evaluate $\nabla_\theta \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$
9:         Compute adapted parameters with gradient descent:
        $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$
10:    Let $p(\theta|\mathcal{D}^{\text{tr}}) = \mathcal{N}(\boldsymbol{\mu}_\theta - \boldsymbol{\gamma}_p \nabla_{\boldsymbol{\mu}_\theta} \mathcal{L}(\boldsymbol{\mu}_\theta, \mathcal{D}^{\text{tr}}), \boldsymbol{\sigma}_\theta^2))$
11:    Compute $\nabla_\Theta \left( \sum_{\mathcal{T}_i} \mathcal{L}(\phi_i, \mathcal{D}^{\text{test}}) \right.$
                $\left. + D_{\text{KL}}(q(\theta|\mathcal{D}^{\text{test}}) \| p(\theta|\mathcal{D}^{\text{tr}})) \right)$
12:    Update $\Theta$ using Adam

Posterior distribution on test data $q(\theta|\mathcal{D}^{\text{test}})$
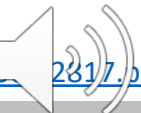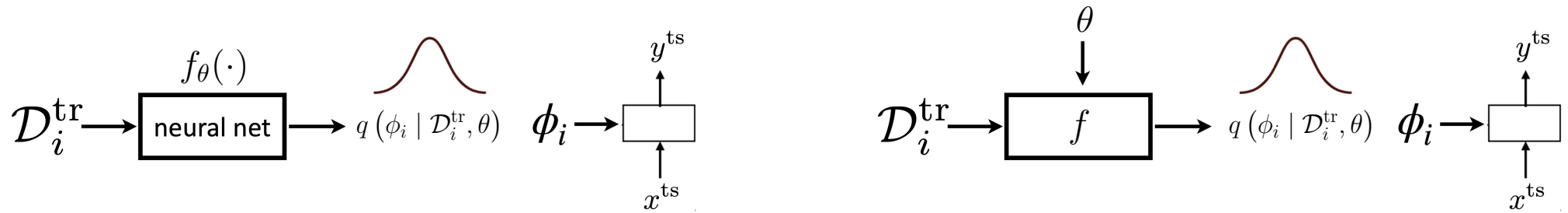
Minimize the gap between two posterior distribution

**Algorithm 2** Meta-testing

**Require:** training data $\mathcal{D}_{\mathcal{T}}^{\text{tr}}$ for new task $\mathcal{T}$
**Require:** learned $\Theta$
1: Sample $\theta$ from the prior $p(\theta|\mathcal{D}^{\text{tr}})$
2: Evaluate $\nabla_\theta \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$
3: Compute adapted parameters with gradient descent:
    $\phi_i = \theta - \alpha \nabla_\theta \mathcal{L}(\theta, \mathcal{D}^{\text{tr}})$

# Amortized variational inference

- **Summary**

- **Latent variable models + variational inference (approximating $p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right)$)**

  approximate likelihood of latent variable model with variational lower bound



$$\max_\theta \; \mathbb{E}_{p(\mathcal{D}_i)}\left[\mathbb{E}_{q(\phi_i|\mathcal{D}_i^{\mathrm{tr}},\theta)}[\ln p(\mathcal{D}_i^{\mathrm{te}}|\phi_i,\theta)] - KL(q(\phi_i|\mathcal{D}_i^{\mathrm{tr}},\theta)\|p(\phi_i|\theta))\right]$$

Pros:

+ can represent non-Gaussian distributions over $y^{\mathrm{ts}}$
+ produces distribution over functions

Cons:

- Can only represent Gaussian distributions $p(\phi_i|\theta)$

Not always restricting: e.g. if $p(y_i^{\mathrm{ts}}|x_i^{\mathrm{ts}}, \phi_i, \theta)$ is also conditioned on $\theta$.

# Amortized variational inference

- **Summary**

- **Latent variable models + variational inference (approximating $p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right)$)**

  approximate likelihood of latent variable model with variational lower bound

  1. $\theta \sim \mathcal{N}(\mu_\theta, \Sigma_\theta)$
  2. $\phi_i \sim p(\phi_i | \theta, x_i^{\text{train}}, y_i^{\text{train}}) \approx \hat{\phi}_i = \theta + \alpha \nabla_\theta \log p(y_i^{\text{train}} | x_i^{\text{train}}, \theta)$

$$\min_{\lambda_\theta} \left( \sum_{\mathcal{T}_i} \mathcal{L}\left(\phi_i, \mathcal{D}^{\text{te}}\right) + KL\left(q\left(\theta \mid \mathcal{D}^{\text{te}}\right) \| p\left(\theta \mid \mathcal{D}^{\text{tr}}\right)\right) \right)$$

**Pros**: Non-Gaussian posterior, simple
at test time, only one model instance.

**Con**: More complex training procedure.

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)
- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks
- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
  - Active-learning evaluation

# Bayesian ensembles

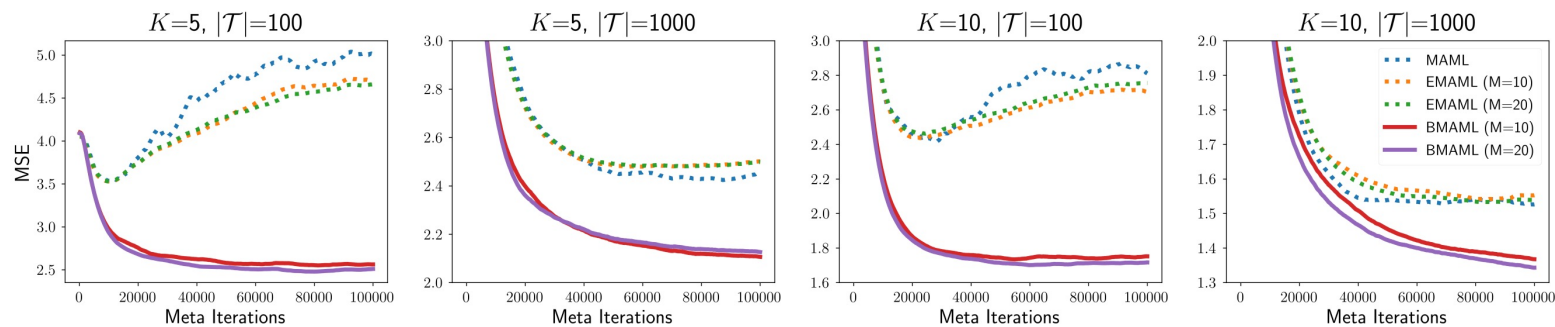- Key idea: train separate models on bootstraps of the data

- Ensemble of MAML (EMAML)

  Train M independent MAML models then take the average

  Does not work as the ensemble members are too similar

- Bayesian Meta-Learning with Chaser Loss (BMAML)

  Use stein variational gradient descent (SVGD) to push the particles away from one another

  Use chaser loss to improve the generalization ability



**Figure 1:** Sinusoidal regression experimental results (meta-testing performance) by varying the number of examples ($K$-shot) given for each task and the number of tasks $|\mathcal{T}|$ used for meta-training.

# Bayesian ensembles

- Notations
  - Meta parameter $\quad \theta \Rightarrow \theta_0$
  - Learner's parameter $\quad \phi \Rightarrow \theta_\tau$

- Stein variational gradient descent (SVGD)

$$\theta_{t+1} \leftarrow \theta_t + \epsilon_t \phi(\theta_t) \quad \text{where} \quad \phi(\theta_t) = \frac{1}{M} \sum_{j=1}^{M} \left[ k(\theta_t^j, \theta_t) \nabla_{\theta_t^j} \log p(\theta_t^j) + \nabla_{\theta_t^j} k(\theta_t^j, \theta_t) \right],$$

keep M models

positive-definite kernel

push the models away from one another

- MAML vs. MAML with SVGD

**Algorithm 1** MAML

Sample a mini-batch of tasks $\mathcal{T}_t$ from $p(\mathcal{T})$
**for** each task $\tau \in \mathcal{T}_t$ **do**
$\quad \theta_\tau \leftarrow \text{GD}_n(\theta_0; \mathcal{D}_\tau^{\text{trn}}, \alpha)$
**end for**
$\theta_0 \leftarrow \theta_0 - \beta \nabla_{\theta_0} \sum_{\tau \in \mathcal{T}_t} \mathcal{L}(\theta_\tau; \mathcal{D}_\tau^{\text{val}})$

**Algorithm 2** Bayesian Fast Adaptation

Sample a mini-batch of tasks $\mathcal{T}_t$ from $p(\mathcal{T})$
**for** each task $\tau \in \mathcal{T}_t$ **do**
$\quad \Theta_\tau(\Theta_0) \leftarrow \text{SVGD}_n(\Theta_0; \mathcal{D}_\tau^{\text{trn}}, \alpha)$
**end for**
$\Theta_0 \leftarrow \Theta_0 - \beta \nabla_{\Theta_0} \sum_{\tau \in \mathcal{T}_t} \mathcal{L}_{\text{BFA}}(\Theta_\tau(\Theta_0); \mathcal{D}_\tau^{\text{val}})$

where $\quad \mathcal{L}_{\text{BFA}}(\Theta_\tau(\Theta_0); \mathcal{D}_\tau^{\text{val}}) = \log \left[ \frac{1}{M} \sum_{m=1}^{M} p(\mathcal{D}_\tau^{\text{val}} | \theta_\tau^m) \right]$

https://proceedings.neurips.cc/paper/2018/file/e1021d43911ca2c1845910d84f40ae__ Paper.pdf

# Bayesian ensembles

- Bayesian MAML with Chaser loss

**Algorithm 3** Bayesian Meta-Learning with Chaser Loss (BMAML)

1: Initialize $\Theta_0$
2: **for** $t = 0, \ldots$ until converge **do**
3:     Sample a mini-batch of tasks $\mathcal{T}_t$ from $p(\mathcal{T})$
4:     **for** each task $\tau \in \mathcal{T}_t$ **do**
5:         Compute chaser $\Theta_\tau^n(\Theta_0) = \text{SVGD}_n(\Theta_0; \mathcal{D}_\tau^{\text{trn}}, \alpha)$
6:         Compute leader $\Theta_\tau^{n+s}(\Theta_0) = \text{SVGD}_s(\Theta_\tau^n(\Theta_0); \mathcal{D}_\tau^{\text{trn}} \cup \mathcal{D}_\tau^{\text{val}}, \alpha)$
7:     **end for**
8:     $\Theta_0 \leftarrow \Theta_0 - \beta \nabla_{\Theta_0} \sum_{\tau \in \mathcal{T}_t} d_s(\Theta_\tau^n(\Theta_0) \parallel \text{stopgrad}(\Theta_\tau^{n+s}(\Theta_0)))$
9: **end for**

- Chaser loss

$$\mathcal{L}_{\text{BMAML}}(\Theta_0) = \sum_{\tau \in \mathcal{T}_t} d_s(\Theta_\tau^n \parallel \Theta_\tau^{n+s}) = \sum_{\tau \in \mathcal{T}_t} \sum_{m=1}^{M} \|\theta_\tau^{n,m} - \theta_\tau^{n+s,m}\|_2^2$$

| $\delta(\phi_i^*)$ | $p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right)$ |
| $\delta(\theta^*)$ | $p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right)$ |

http://cs330.stanford.edu/fall2020/slides/cs330_bayesian_meta_learning_2020.pdf
https://proceedings.neurips.cc/paper/2018/file/e1021d43911ca2c1845910d84f40ae_Paper.pdf

# Bayesian ensembles

- Bayesian MAML with Chaser loss

**Algorithm 3** Bayesian Meta-Learning with Chaser Loss (BMAML)

1: Initialize $\Theta_0$
2: **for** $t = 0, \ldots$ until converge **do**
3:   Sample a mini-batch of tasks $\mathcal{T}_t$ from $p(\mathcal{T})$
4:   **for** each task $\tau \in \mathcal{T}_t$ **do**
5:     Compute chaser $\Theta_\tau^n(\Theta_0) = \text{SVGD}_n(\Theta_0; \mathcal{D}_\tau^{\text{trn}}, \alpha)$
6:     Compute leader $\Theta_\tau^{n+s}(\Theta_0) = \text{SVGD}_s(\Theta_\tau^n(\Theta_0); \mathcal{D}_\tau^{\text{trn}} \cup \mathcal{D}_\tau^{\text{val}}, \alpha)$
7:   **end for**
8:   $\Theta_0 \leftarrow \Theta_0 - \beta \nabla_{\Theta_0} \sum_{\tau \in \mathcal{T}_t} d_s(\Theta_\tau^n(\Theta_0) \parallel \text{stopgrad}(\Theta_\tau^{n+s}(\Theta_0)))$
9: **end for**

- Chaser loss

$$\mathcal{L}_{\text{BMAML}}(\Theta_0) = \sum_{\tau \in \mathcal{T}_t} d_s(\Theta_\tau^n \parallel \Theta_\tau^{n+s}) = \sum_{\tau \in \mathcal{T}_t} \sum_{m=1}^{M} \|\theta_\tau^{n,m} - \theta_\tau^{n+s,m}\|_2^2$$

**Pros**: Simple, tends to work well, non-Gaussian distributions.

**Con**: Need to maintain M model instances. (or do gradient-based inference on **last layer only**)

$\delta(\phi_i^*)$ ✗   $p(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*)$ ✓

$\delta(\theta^*)$ ✗   $p(\theta \mid \mathcal{D}_{\text{meta-train}})$ ✓

http://cs330.stanford.edu/fall2020/slides/cs330_bayesian_meta_learning_2020.pdf
https://proceedings.neurips.cc/paper/2018/file/e1021d43911ca2c1845910d84f40ae_Paper.pdf

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)
- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks
- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
  - Active-learning evaluation

# Bayesian neural networks

- Key idea: explicit distribution over the space of network parameters

- Monte Carlo dropout in neural networks can be used to perform variational inference to make it Bayesian

model parameter is sampled by dropping different neurons

$q(\boldsymbol{\omega})$ :

$$\mathbf{W}_i = \mathbf{M}_i \cdot \mathrm{diag}([\mathbf{z}_{i,j}]_{j=1}^{K_i})$$

$$\mathbf{z}_{i,j} \sim \mathrm{Bernoulli}(p_i) \text{ for } i = 1, ..., L, \ j = 1, ..., K_{i-1}$$

Easiest way to get a distribution approximating the

posterior distribution

$$q(\boldsymbol{\omega}) \overset{\text{Approx.}}{\rightarrow} p(\boldsymbol{\omega}|\mathcal{D}^{\mathrm{tr}})$$

$$\hat{w} = w_i \quad \textbf{or} \quad \hat{w} = w_j$$

# Bayesian neural networks

Combine the SGD with Bayesian neural networks (AGILE)

$$q(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta) \rightarrow p(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta)$$

1. Initial learner's parameter with $\theta$

2. Calculate the inner loop loss on $\mathcal{D}_i^{\mathrm{tr}}$ using model with dropout parametrized by $\phi_i$

3. Update parameter $\phi_i$ using SGD

4. Go back to step 2 for more gradient descent steps

Loss function:

- Inner loop loss for optimizing $\phi_i$ : $\lambda_i = \arg\max_{\lambda_i} \log p\left(\mathcal{D}_i^{\mathrm{tr}} \mid \phi_i\right) + \log p\left(\phi_i \mid \theta\right)$ , $\lambda_i = \mathrm{M}_i$ is a set of complete parameter (no dropout)

- Outer loop loss for optimizing $\theta$ :

$$\theta^* = \underset{\theta \in \Theta}{\arg\max} \prod_{i=1}^{M} p(\mathcal{D}_i^{\mathrm{te}}|\mathcal{D}_i^{\mathrm{tr}}, \theta)$$

$$= \underset{\theta \in \Theta}{\arg\max} \prod_{i=1}^{M} \left( \int p(\mathcal{D}_i^{\mathrm{te}}|\phi_i)p(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta)d\phi_i \right)$$

$$\approx \underset{\theta \in \Theta}{\arg\max} \prod_{i=1}^{M} \left( \frac{1}{T} \sum_{t=1}^{t} p(y_i^{\mathrm{te}}|x_i^{\mathrm{te}}, \phi_i^t) \right), \quad \text{where } \phi_i^t \sim q(\phi_i|\mathcal{D}_i^{\mathrm{tr}}, \theta).$$

https://arxiv.org/pdf/200~~~~~~~df.

# Bayesian neural networks

Combine the SGD with Bayesian neural networks (AGILE)

- Dropout as well during test

$$p(\mathbf{y}^{\text{te}}|\mathbf{x}^{\text{te}}, \mathcal{D}^{\text{tr}}, \theta) = \int p(\mathbf{y}^{\text{te}}|\mathbf{x}^{\text{te}}, \phi)q(\phi|\mathcal{D}^{\text{tr}}, \theta)d\phi$$

$$\approx \frac{1}{T}\sum_{t=1}^{t} p(y^{\text{te}}|x^{\text{te}}, \phi^t), \quad \text{where } \phi^t \sim q(\phi|\mathcal{D}^{\text{tr}}, \theta)$$

Not the complete
parameter set $\mathbb{M}$

**Pros:** Simple, only one model instance

**Cons:** Can only model Gaussian distribution (Bayesian neural network),
need to finetune the hyperparameter dropout rate

$$\delta(\phi_i^*) \quad \textsf{X} \qquad p\left(\phi_i \mid \mathcal{D}_i^{\text{tr}}, \theta^*\right) \qquad \checkmark$$
$$\delta(\theta^*) \quad \checkmark \qquad p\left(\theta \mid \mathcal{D}_{\text{meta-train}}\right) \qquad \textsf{X}$$

# Bayesian meta-learning approaches summary

- **Latent variable models + variational inference**

    approximate likelihood of latent variable model with variational lower bound

    **Approximating** $\qquad p\left(\phi_i \mid \mathcal{D}_i^{\mathrm{tr}}, \theta^*\right) \qquad\qquad\qquad\qquad p\left(\theta \mid \mathcal{D}_{\mathrm{meta\text{-}train}}\right)$

    Pros:

    \+ can represent non-Gaussian distributions over $y^{\mathrm{ts}}$
    \+ produces distribution over functions

    Cons:

    \- Can only represent Gaussian distributions $p(\phi_i \mid \theta)$

    Pros: Non-Gaussian posterior, simple
    at test time, only one model instance.

    Con: More complex training procedure.

- **Bayesian ensembles**

    particle-based representation: train separate models on bootstraps of the data

    Pros: Simple, tends to work well,
    non-Gaussian distributions.

    Con: Need to maintain M model instances.
    (or do gradient-based inference on **last layer only**)

- **Bayesian neural networks**

    explicit distribution over the space of network parameters

    Pros: Simple, only one model instance

    Cons: Can only model Gaussian distribution (Bayesian neural network),
    need to finetune the hyperparameter dropout rate

    http://cs330.stanford.edu/fall2020/slides/cs330_bayesian_meta_learning_2020.pdf

# Outline

- Introduction
  - Why Bayesian meta-learning?
  - The evidence lower bound (ELBO)

- Bayesian meta-learning approaches based on
  - Amortized variational inference
    - Black-box
    - Optimization
  - Bayesian ensembles
  - Bayesian neural networks

- Bayesian meta-learning evaluation
  - Qualitative visualization
  - Quantitative evaluation
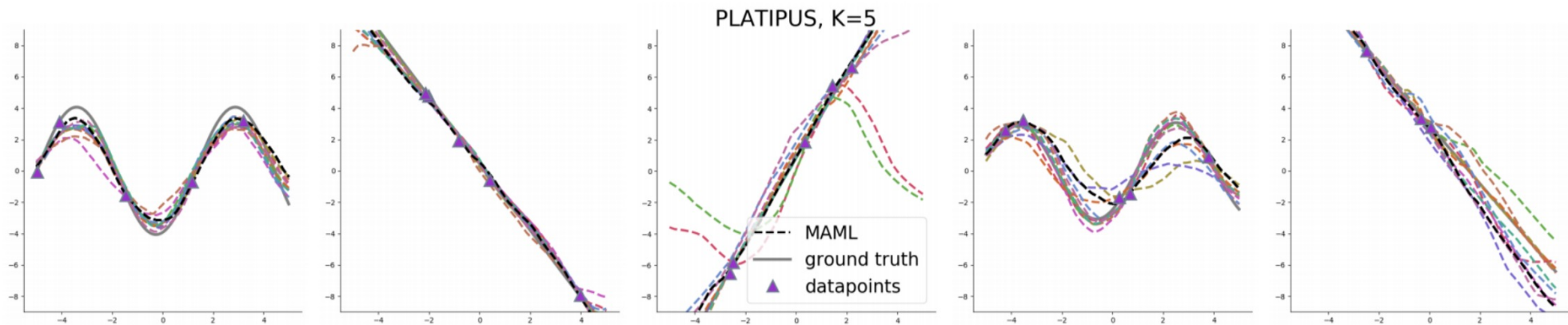  - Active-learning evaluation
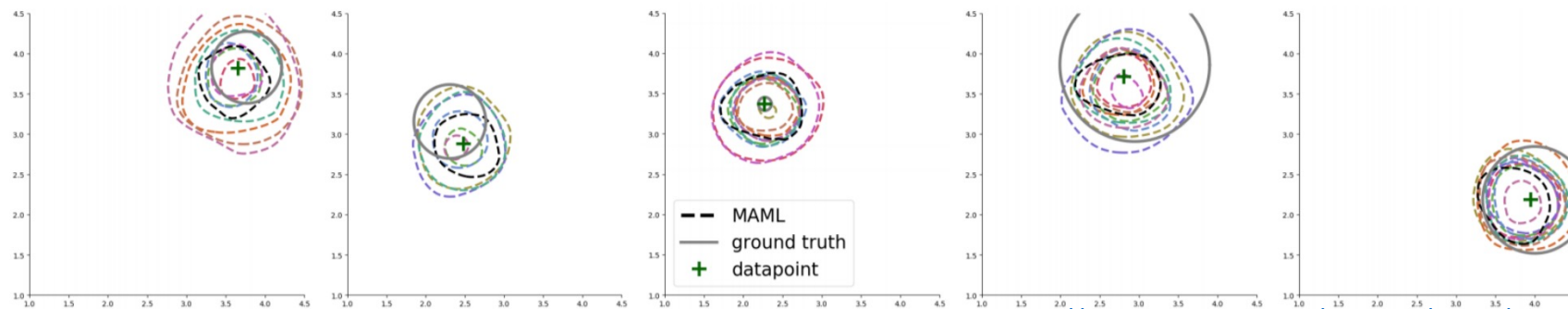
# Qualitative visualization

## Qualitative Evaluation on Toy Problems with Ambiguity

**(Finn\*, Xu\*, Levine, NeurIPS '18)**

Ambiguous regression:



Ambiguous classification:

# Quantitative evaluation

## Evaluation on Ambiguous Generation Tasks

(Gordon et al., ICLR '19)



| Model | MSE | SSIM |
|---|---|---|
| C-VAE 1-shot | 0.0269 | 0.5705 |
| VERSA 1-shot | 0.0108 | 0.7893 |
| VERSA 5-shot | 0.0069 | 0.8483 |

**Table 2:** View reconstruction test results.

# Quantitative evaluation

## Accuracy, Mode Coverage, & Likelihood on Ambiguous Tasks

(Finn*, Xu*, Levine, NeurIPS '18)



| Ambiguous celebA (5-shot) | | | |
|---|---|---|---|
| | Accuracy | Coverage (max=3) | Average NLL |
| MAML | **89.00 ± 1.78**% | 1.00 ± 0.0 | 0.73 ± 0.06 |
| MAML + noise | 84.3 ± 1.60 % | 1.89 ± 0.04 | 0.68 ± 0.05 |
| **PLATIPUS (ours) (KL weight = 0.05)** | **88.34 ± 1.06** % | 1.59 ± 0.03 | 0.67± 0.05 |
| **PLATIPUS (ours) (KL weight = 0.15)** | **87.8 ± 1.03** % | **1.94 ± 0.04** | **0.56 ± 0.04** |

# Quantitative evaluation

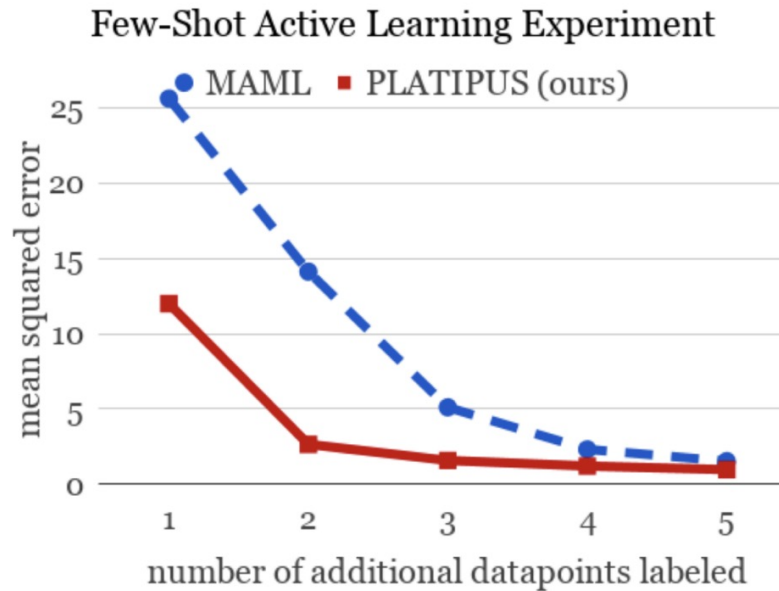## Reliability Diagrams & Accuracy

(Ravi & Beatson, ICLR '19)
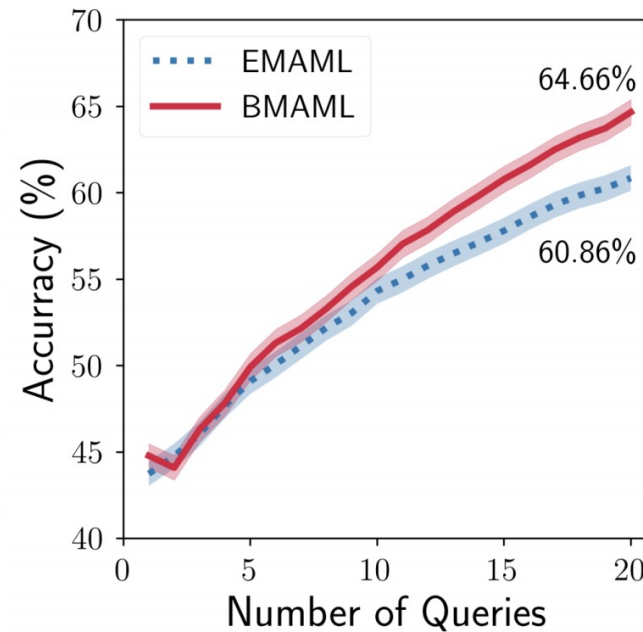


*mini*ImageNet: 1-shot, 5-class

| *mini*ImageNet | 1-shot, 5-class |
|---|---|
| MAML (ours) | $47.0 \pm 0.59$ |
| Prob. MAML (ours) | $47.8 \pm 0.61$ |
| Our Model | $45.0 \pm 0.60$ |

# Active-learning evaluation

**Finn*, Xu*, Levine, NeurIPS '18**
Sinusoid Regression



Few-Shot Active Learning Experiment

**Kim et al. NeurIPS '18**
MiniImageNet



Both experiments:
- Sequentially choose datapoint with
  **maximum predictive entropy** to be labeled
- or choose datapoint at random (MAML)

# Take away

- Uncertainty is important when study meta-learning as few annotated samples are provided for each tasks

- Uncertainty can exist in either meta parameter or learner's parameter or both

- There are several tools to make the meta learning algorithm Bayesian:
  - Latent variable models + variational inference
  - Bayesian ensembles
  - Bayesian neural networks

- Learn different ways to evaluate the Bayesian meta learning algorithm