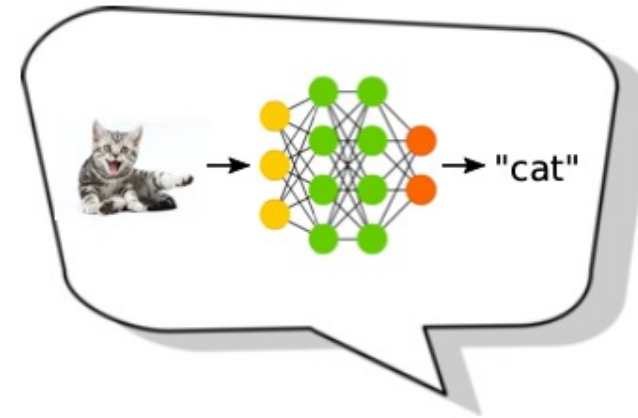
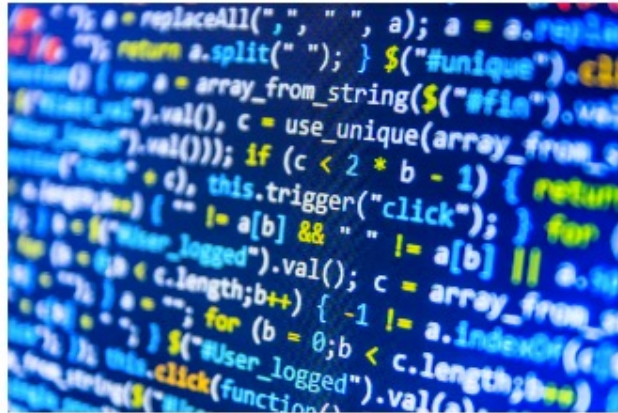
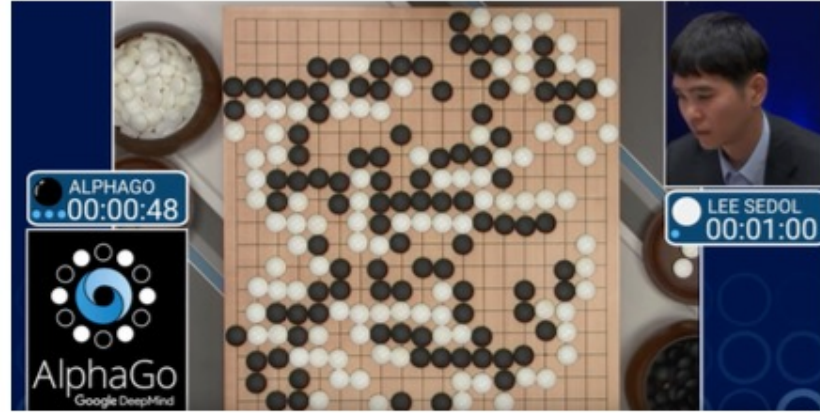


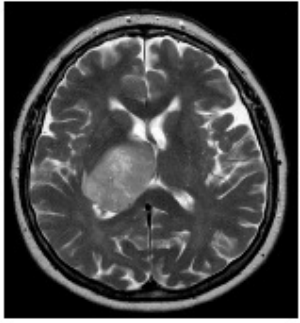
# An Introduction to Bayesian Neural Networks

Yingzhen Li

[yingzhen.li@imperial.ac.uk](mailto:yingzhen.li@imperial.ac.uk)



**Deep Learning**



brain tumor  
type "A"

are you sure? why?



Do you know what  
you don't know?  
How confident are you?



# Bayesian Inference

$$\pi(\theta) = p(\theta | \text{data})$$

$$P(\theta | \text{data}) = \frac{P(\theta)P(\text{data} | \theta)}{P(\text{data})}$$

- $P(\theta)$ : prior distribution
- $P(\text{data} | \theta)$ : likelihood of  $\theta$  given  $\text{data}$
- $P(\theta | \text{data})$ : posterior distribution of  $\theta$  given  $\text{data}$
- $P(\text{data})$ : marginal likelihood/model evidence

$$P(\text{data}) = \int P(\theta)P(\text{data} | \theta)$$



Image courtesy of Sebastian Nowozin

Re-use of the image for any other purpose is not allowed

# Bayesian Inference

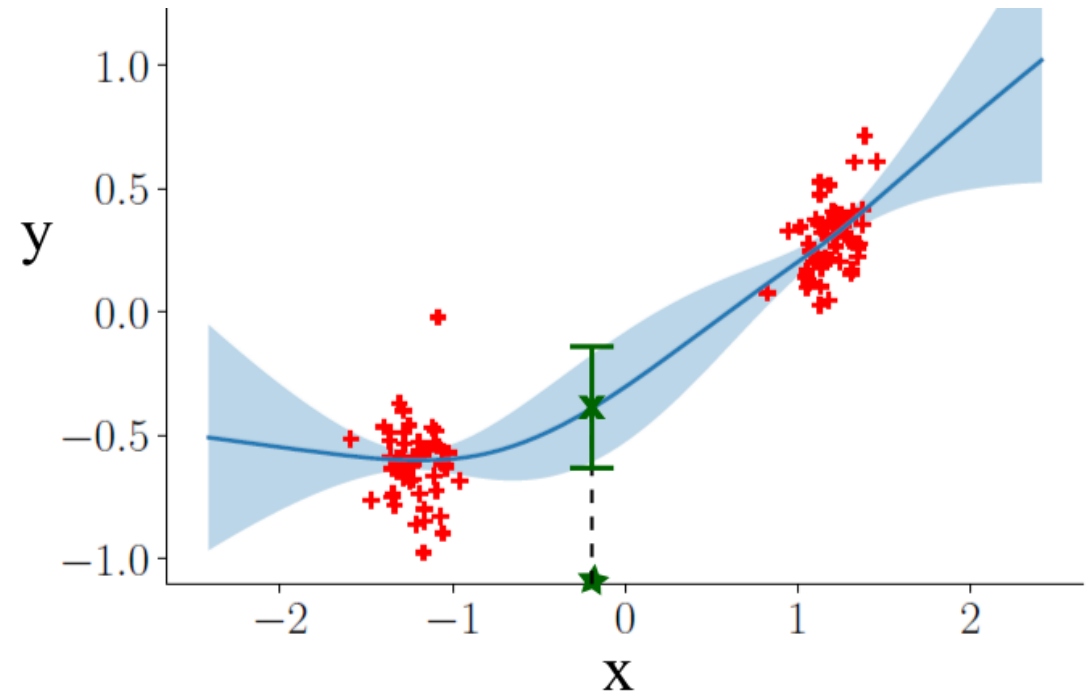
- The central equation for Bayesian inference:

$$\int F(\theta)p(\theta|D)d\theta$$

“What is the prediction distribution of the **test output** given a **test input**?”

$$F(\theta) = p(y|x, \theta),$$

$D$  = observed datapoints

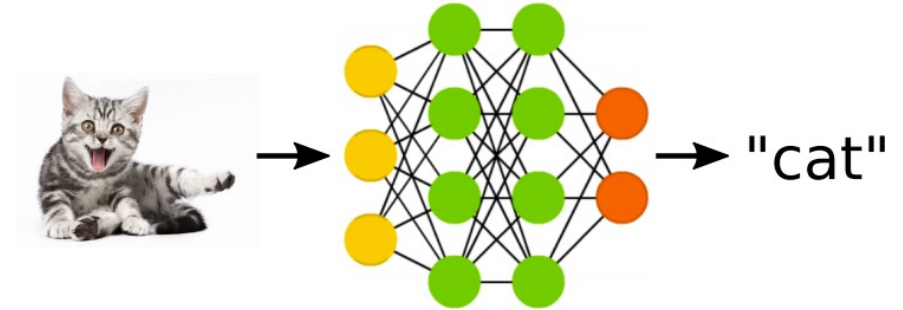


# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$ : input image;  $y$ : output label
- Build a neural network with parameters  $\theta$ :

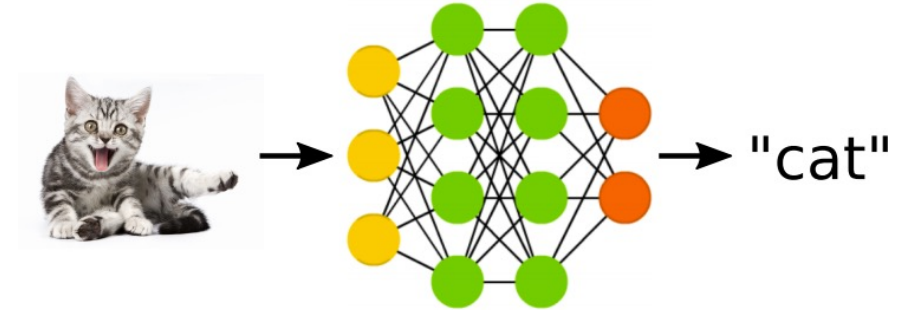
$$p(y|x, \theta) = \text{softmax}(f_{\theta}(x))$$



# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$ : input image;  $y$ : output label
- Build a neural network with parameters  $\theta$ :  
$$p(y|x, \theta) = \text{softmax}(f_{\theta}(x))$$



A typical neural network (with non-linearity  $g(\cdot)$ ):

$$f_{\theta}(x) = W^L g(W^{L-1} g(\dots g(W^1 x + b^1)) + b^{L-1}) + b^L,$$

$$h^l = g(W^l h^{l-1} + b^l), h^1 = g(W^1 x + b^1).$$

Neural network parameters:  $\theta = \{W^l, b^l\}_{l=1}^L$

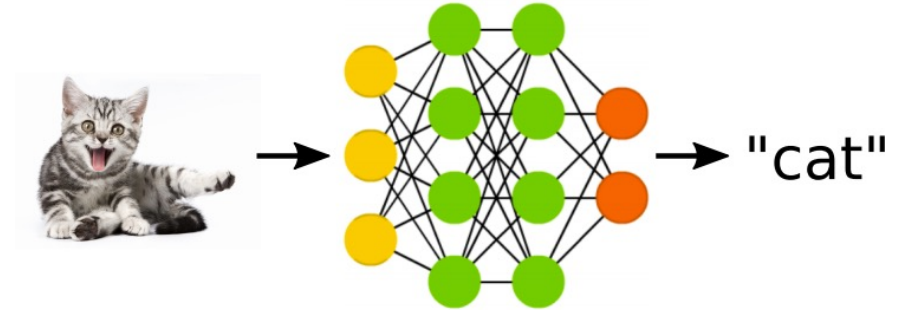


# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$ : input image;  $y$ : output label
- Build a neural network with parameters  $\theta$ :

$$p(y|x, \theta) = \text{softmax}(f_{\theta}(x))$$



Typical deep learning solution:

- Optimize  $\theta$  to obtain a point estimates (MLE):

$$\theta^* = \operatorname{argmax} \log p(D | \theta),$$
$$\log p(D | \theta) = \sum_{n=1}^N \log p(y_n | x_n, \theta), D = \{(x_n, y_n)\}_{n=1}^N$$

- Prediction: using  $p(y^* | x^*, \theta^*)$

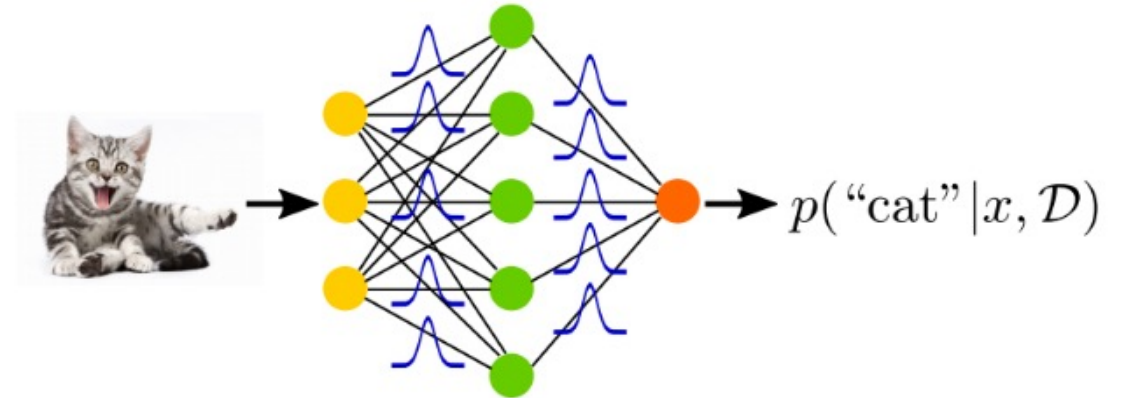


# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$ : input image;  $y$ : output label
- Build a neural network with parameters  $\theta$ :

$$p(y|x, \theta) = \text{softmax}(f_{\theta}(x))$$



## Bayesian solution:

- Put a prior  $p(\theta)$  on network parameters  $\theta$ , e.g. Gaussian prior

$$p(\theta) = N(\theta; 0, \sigma^2 I)$$

- Compute the posterior distribution  $p(\theta | D)$ :

$$p(\theta | D) \propto p(D | \theta) p(\theta)$$

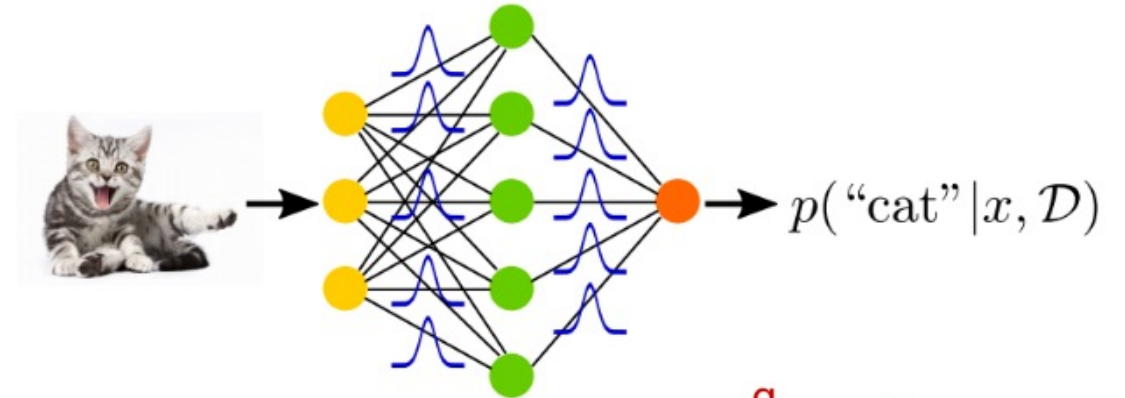
- Bayesian predictive inference:

$$p(y^* | x^*, D) = E_{p(\theta | D)}[p(y^* | x^*, \theta)]$$

# Bayesian Neural Network (BNN) 101

Classifying different types of animals:

- $x$ : input image;  $y$ : output label
- Build a neural network with parameters  $\theta$ :  
$$p(y|x, \theta) = \text{softmax}(f_{\theta}(x))$$



**Approximate (Bayesian) inference solution:**

- Exact posterior intractable, use approximate posterior:

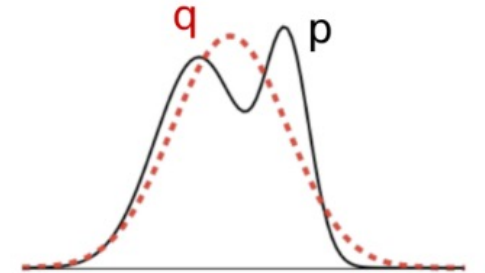
$$q(\theta) \approx p(\theta | D)$$

- Approximate Bayesian predictive inference:

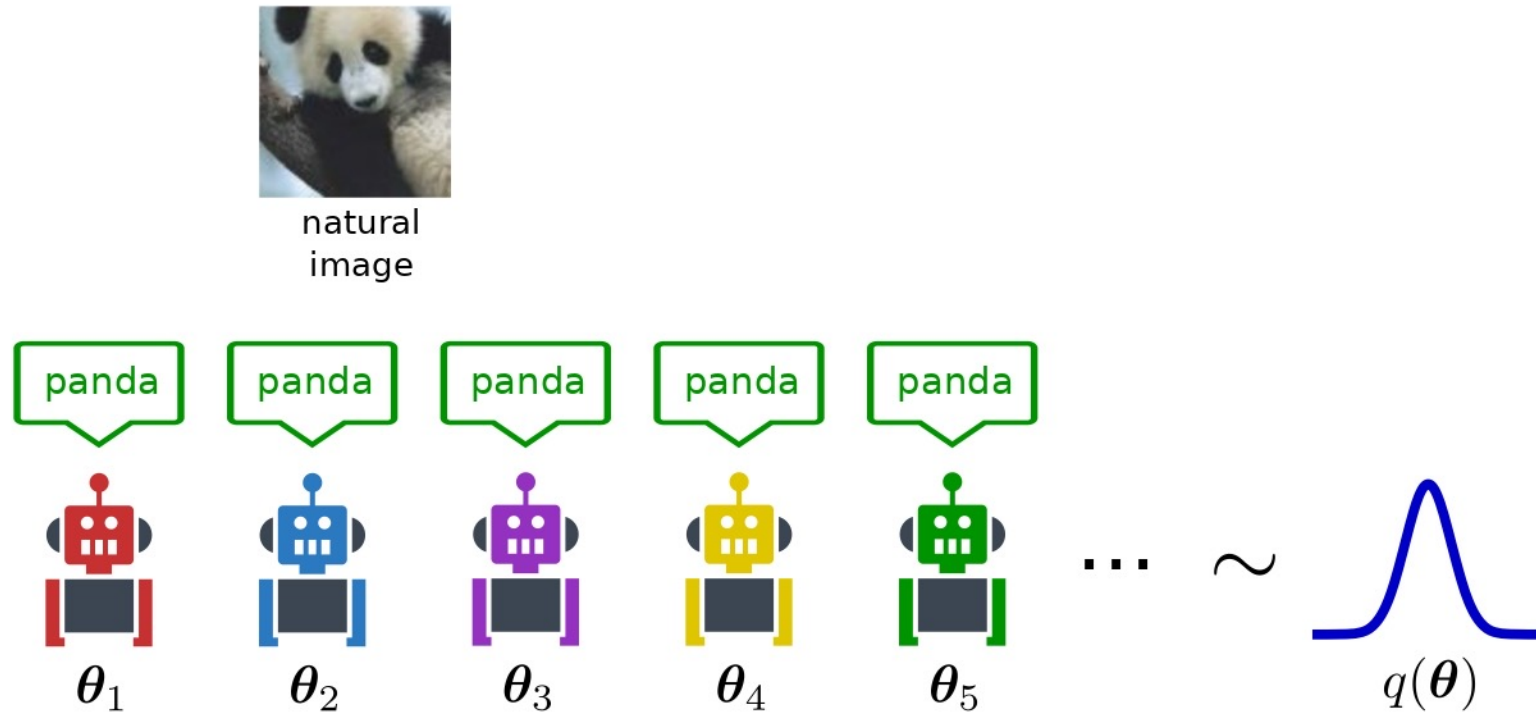
$$p(y^* | x^*, D) \approx E_{q(\theta)}[p(y^* | x^*, \theta)]$$

- Monte Carlo approximation:

$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$

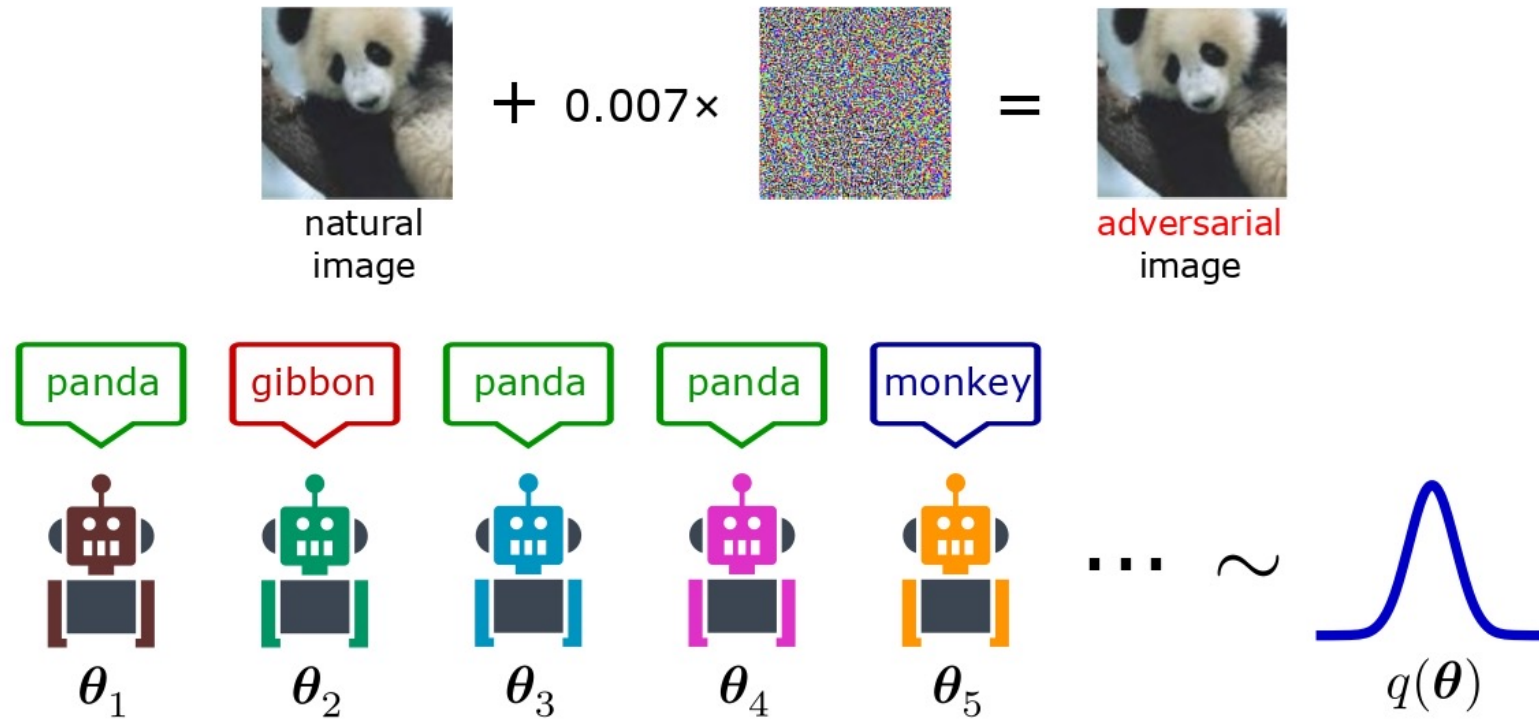


# Bayesian Neural Network (BNN) 101



Prediction on in-distribution data:  
ensemble over networks, using weights sampled from  $q(\theta)$

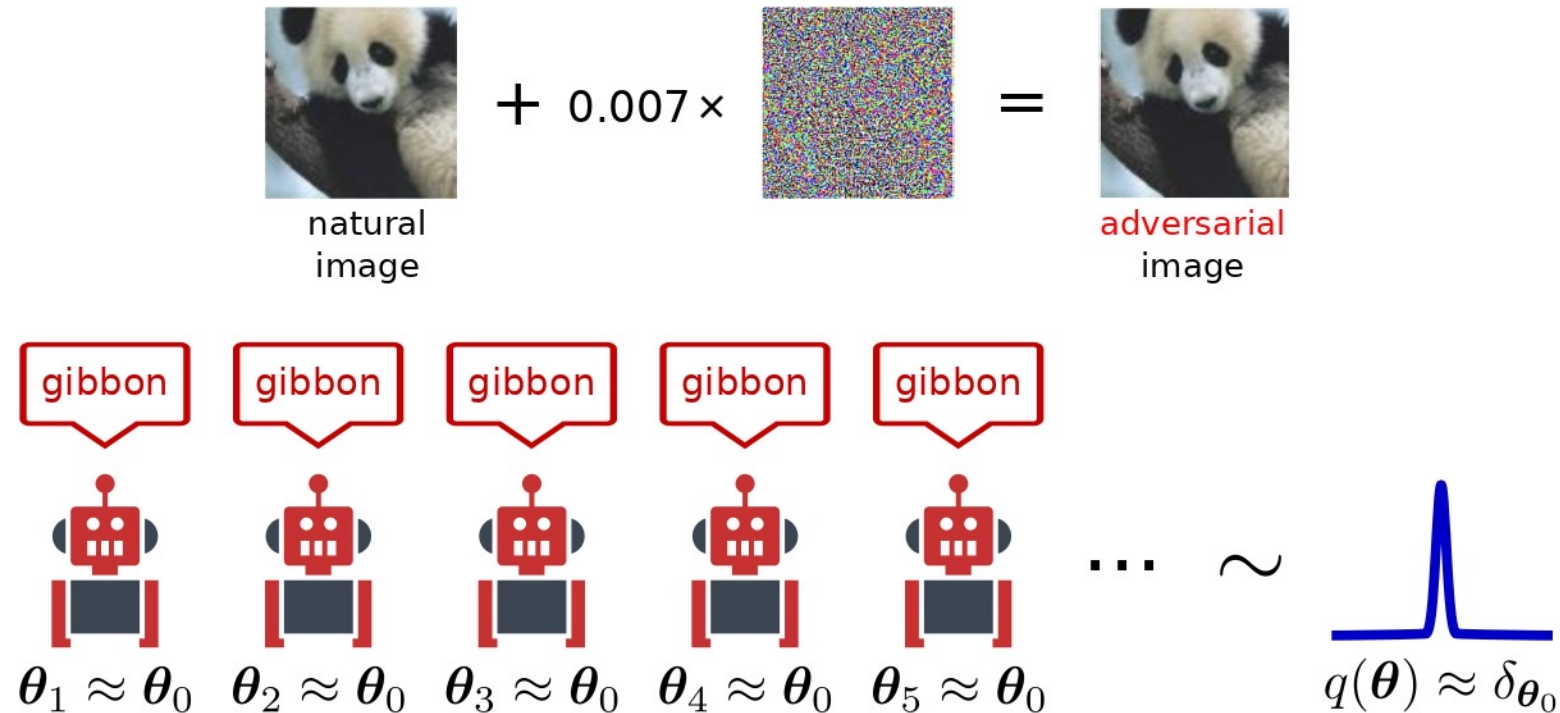
# Bayesian Neural Network (BNN) 101



Prediction on OOD/noisy/adversarial data:

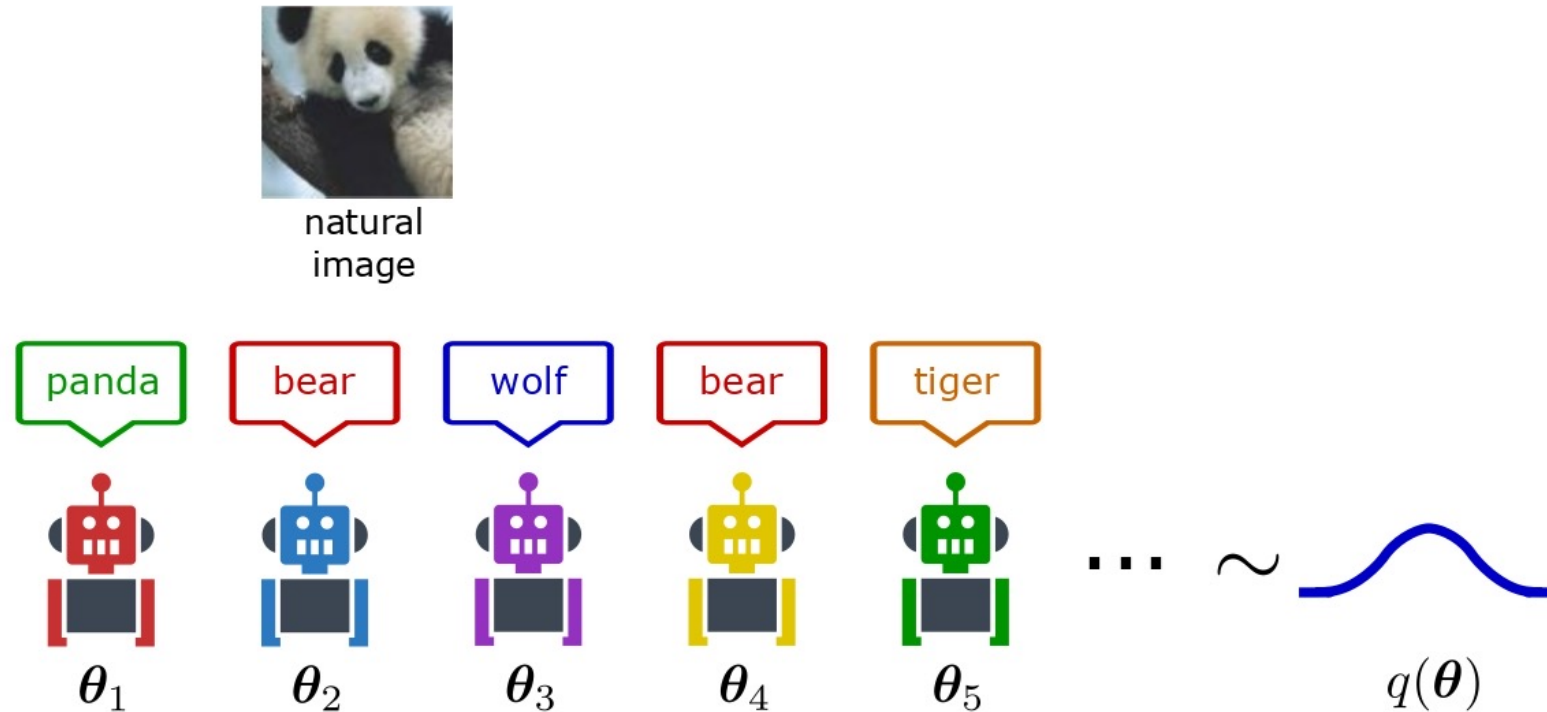
Disagreement (i.e. uncertainty) exists over networks sampled from  $q(\theta)$

# Bayesian Neural Network (BNN) 101



Prediction on OOD/noisy/adversarial data **when  $q(\theta)$  is over-confident:**  
Return **confidently wrong answers** (close to point estimate)

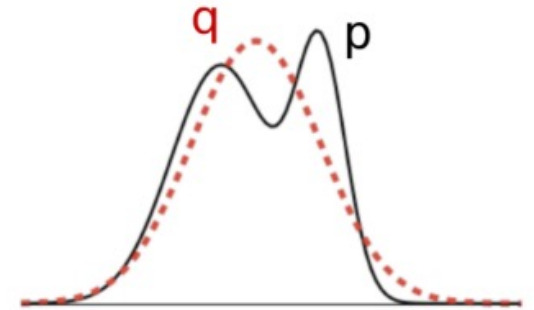
# Bayesian Neural Network (BNN) 101



Prediction on in-distribution data **when  $q(\theta)$  is under-confident:**  
**Low accuracy** in prediction tasks (less desirable)

# Approximate Inference in BNNs

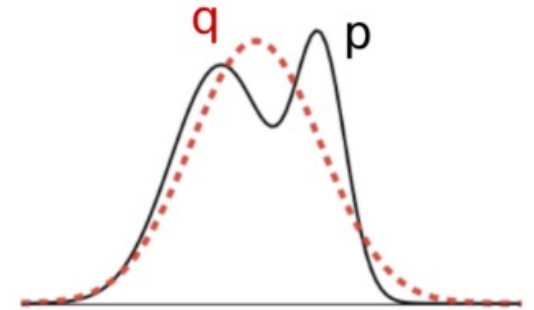
- Key steps of approximate inference in BNNs
  1. Construct the  $q(\theta) \approx p(\theta \mid D)$  distribution
    - Simple distributions: e.g. Mean-field Gaussian
    - Structured approximations, e.g. low-rank Gaussians
    - Others (non-Gaussian)





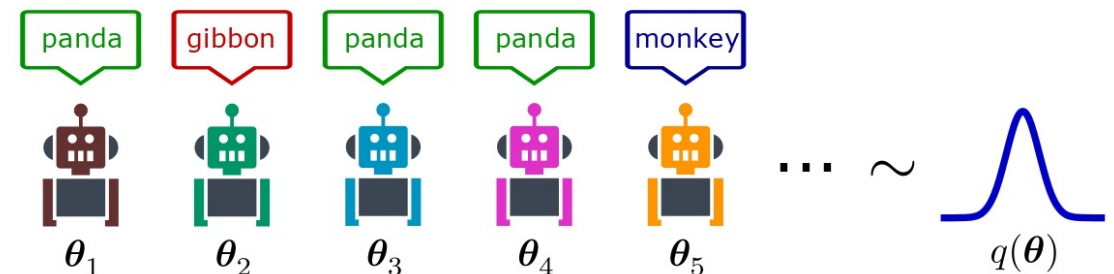
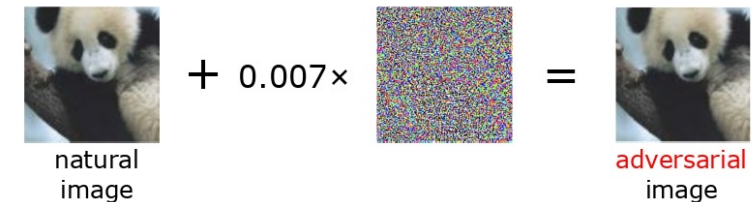
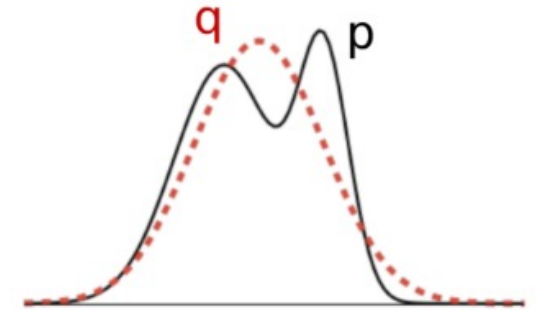
# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
  1. Construct the  $q(\theta) \approx p(\theta \mid D)$  distribution
    - Simple distributions: e.g. Mean-field Gaussian
    - Structured approximations, e.g. low-rank Gaussians
    - Others (non-Gaussian)
  2. Fit the  $q(\theta)$  distribution
    - E.g. with variational inference



# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
  1. Construct the  $q(\theta) \approx p(\theta \mid D)$  distribution
    - Simple distributions: e.g. Mean-field Gaussian
    - Structured approximations, e.g. low-rank Gaussians
    - Others (non-Gaussian)
  2. Fit the  $q(\theta)$  distribution
    - E.g. with variational inference
  3. Compute prediction with Monte Carlo approximations



# Today's agenda

- Lecture on Basics: MFVI for BNNs
- Hands-on tutorial on BNNs
  - i.e., programming exercises
  - Also some case studies

The slide features decorative curved lines in the top corners. On the top left, a light blue curve starts from the edge and curves towards the center. On the top right, a light green curve starts from the edge and curves towards the center. Both curves have a slight gradient and a soft shadow effect.

# Part I: Basics

- **Variational inference**
- **Bayes-by-backprop**

# Bayesian Inference

$$P(\theta | D) = \frac{P(\theta)P(D | \theta)}{P(D)}$$

- $P(\theta)$ : prior
- $P(D | \theta)$ : likelihood
- $P(\theta | D)$ : posterior
- $P(D)$ : marginal



# Variational Inference (VI)

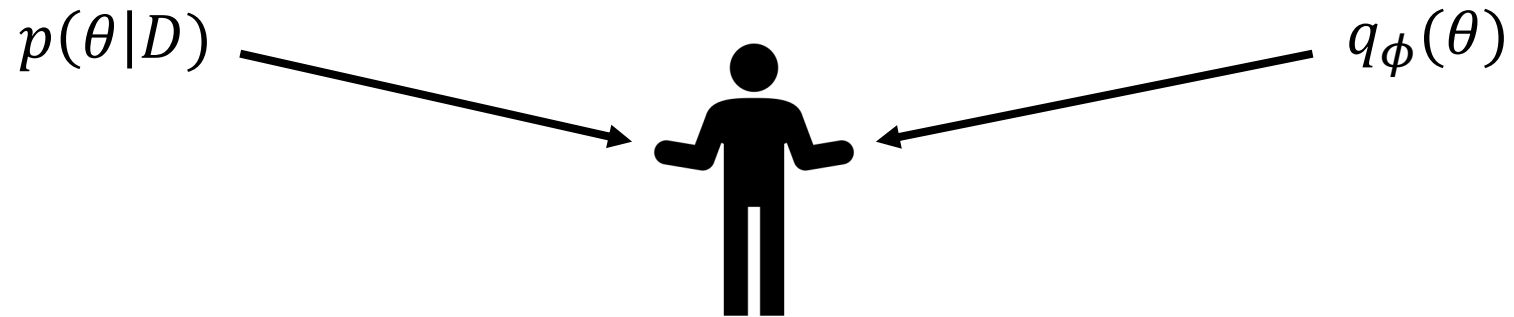
The posterior

$$p(\theta|D) = p(D|\theta)p(\theta)/p(D)$$

The variational distribution

$$q_{\phi}(\theta)$$

# Inference as Optimization



Kullback-Leibler (KL) divergence



# Kullback-Leibler Divergence

$$KL[q(\theta)||p(\theta)] = \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta = E_{q(\theta)}[\log \frac{q(\theta)}{p(\theta)}]$$

- When  $p = q$ , KL is 0
- Otherwise,  $KL > 0$
- It measures how similar are these two distributions

# Let's Derive the Objective of VI

- Minimize  $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -E_{q(\theta)} \left[ \log \frac{p(\theta|D)}{q(\theta)} \right]$$

# Let's Derive the Objective of VI

- Minimize  $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -E_{q(\theta)} \left[ \log \frac{p(\theta|D)}{q(\theta)} \right]$$

$$= -E_{q(\theta)} \left[ \log \frac{p(\theta,D)}{p(D)q(\theta)} \right] = -E_{q(\theta)} \left[ \log \frac{p(\theta,D)}{q(\theta)} - \log p(D) \right]$$

# Let's Derive the Objective of VI

- Minimize  $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = -E_{q(\theta)} \left[ \log \frac{p(\theta|D)}{q(\theta)} \right]$$

$$= -E_{q(\theta)} \left[ \log \frac{p(\theta,D)}{p(D)q(\theta)} \right] = -E_{q(\theta)} \left[ \log \frac{p(\theta,D)}{q(\theta)} - \log p(D) \right]$$

$$= \boxed{\log p(D)} - E_{q(\theta)} \left[ \log \frac{p(\theta,D)}{q(\theta)} \right]$$

Model Evidence

# Let's Derive the Objective of VI

Minimize  $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = \log p(D) - E_{q(\theta)} \left[ \log \frac{p(\theta, D)}{q(\theta)} \right]$$

Maximize  $E_{q(\theta)} \left[ \log \frac{p(\theta, D)}{q(\theta)} \right]$

# Let's Derive the Objective of VI

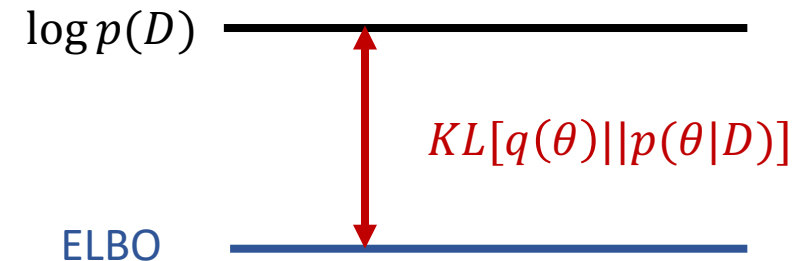
Minimize  $KL[q(\theta)||p(\theta|D)]$

$$KL[q(\theta)||p(\theta|D)] = \boxed{\log p(D)} - E_{q(\theta)} \left[ \log \frac{p(\theta, D)}{q(\theta)} \right]$$

Model Evidence

Maximize  $L = E_{q(\theta)} \left[ \log \frac{p(\theta, D)}{q(\theta)} \right]$

Evidence Lower Bound (ELBO)



“Model Evidence = ELBO + KL”

# Variational Inference (VI)

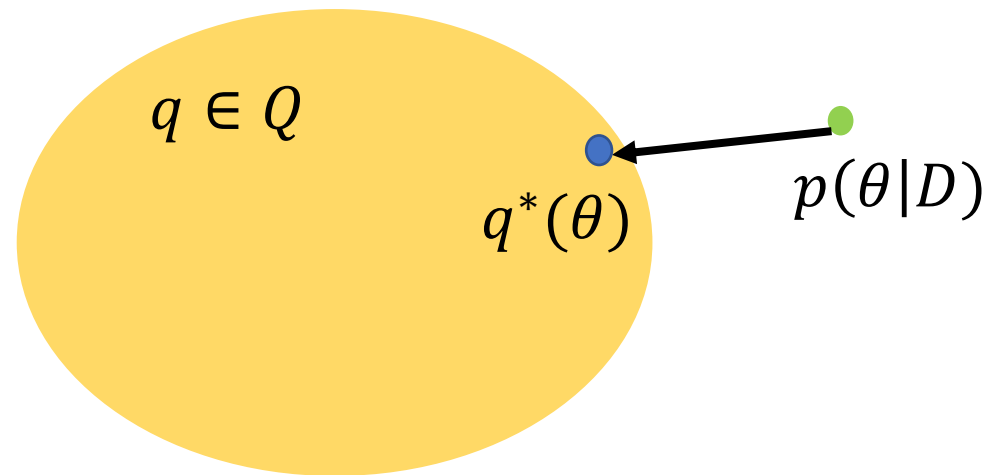
The posterior

$$p(\theta|D) = p(D|\theta)p(\theta)/p(D)$$

The variational distribution

$$q_{\phi}(\theta)$$

$$L = E_{q_{\phi}(\theta)} \left[ \log \frac{p(D, \theta)}{q_{\phi}(\theta)} \right] = \log p(D) - KL[ q_{\phi}(\theta) || p(\theta) ]$$

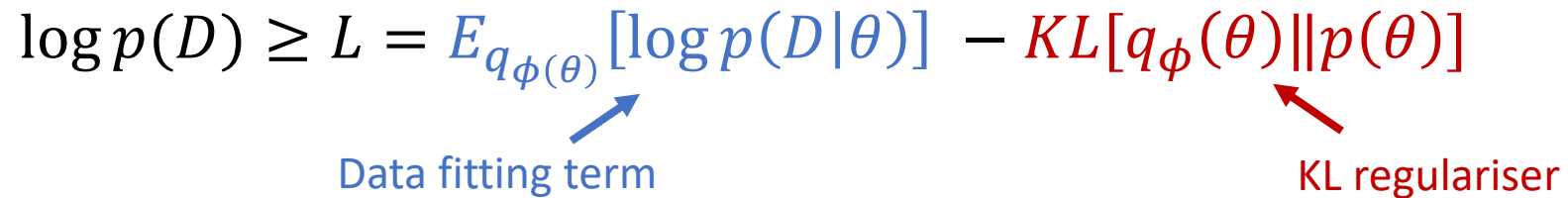




# Variational Inference (VI)

- Rewriting the ELBO:

$$\log p(D) \geq L = E_{q_{\phi}(\theta)}[\log p(D|\theta)] - KL[q_{\phi}(\theta) \| p(\theta)]$$

  
Data fitting term                      KL regulariser

## (Negative) Data fitting term:

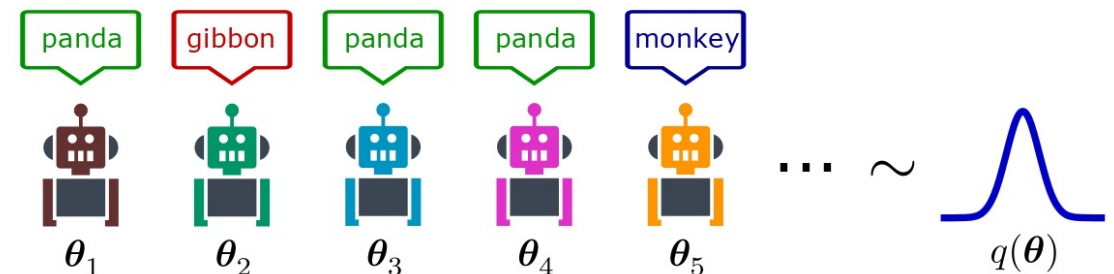
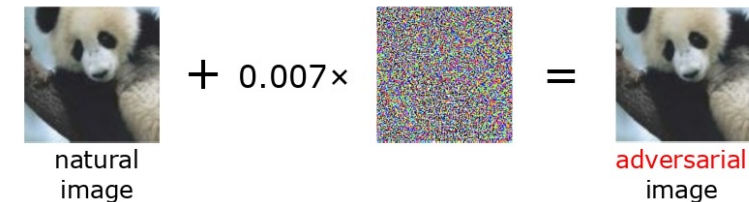
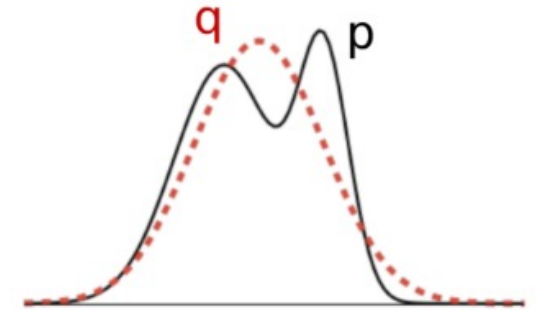
- Like the usual DL loss you'll use for training neural networks
- ...except that now the network's weights are sampled from  $q$

## KL regulariser:

- Make the  $q$  distribution closer to the prior
- Regularises the approximate posterior, especially when using e.g., Gaussian prior

# Approximate Inference in BNNs

- Key steps of approximate inference in BNNs
  1. Construct the  $q(\theta) \approx p(\theta \mid D)$  distribution
    - Simple distributions: e.g. Mean-field Gaussian
    - Structured approximations, e.g. low-rank Gaussians
    - Others (non-Gaussian)
  2. Fit the  $q(\theta)$  distribution
    - E.g. with variational inference
  3. Compute prediction with Monte Carlo approximations



# Approximate Inference in BNNs

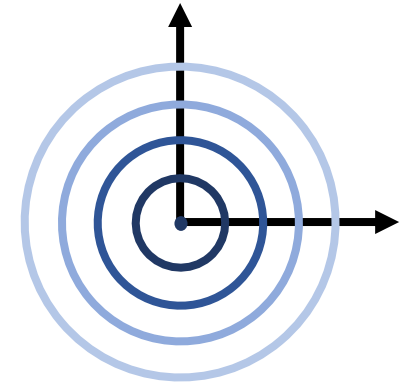
- Step 1: construct the  $q(\theta) \approx p(\theta \mid D)$  distribution
  - Example: Mean-field Gaussian distribution:

$$q(\theta) = \prod_{l=1}^L q(W^l) q(b^l)$$

$$q(W_l) = \prod_{ij} q(W_{ij}^l), \quad q(W_{ij}^l) = N(W_{ij}^l; M_{ij}^l, V_{ij}^l)$$

$$q(b^l) = \prod_i q(b_i^l), \quad q(b_i^l) = N(b_i^l; m_i^l, v_i^l)$$

- Variational parameters:  $\phi = \{M_{ij}^l, \log V_{ij}^l, m_i^l, \log v_i^l\}_{l=1}^L$



# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:
  - Variational inference:  $\phi^* = \operatorname{argmax} L(\phi)$ 
$$L(\phi) = E_{q_\phi(\theta)}[\log p(D \mid \theta)] - KL[q_\phi(\theta) \parallel p(\theta)]$$

# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:

- Variational inference:  $\phi^* = \operatorname{argmax} L(\phi)$

$$L(\phi) = E_{q_\phi(\theta)}[\log p(D \mid \theta)] - KL[q_\phi(\theta) \parallel p(\theta)]$$

- First scalable technique: **Stochastic optimization**

- i.i.d. assumption of data:  $\log p(D \mid \theta) = \sum_{n=1}^N \log p(y_n \mid x_n, \theta)$
- Enable mini-batch training with  $\{(x_m, y_m)\} \sim D^M$ :

$$L(\phi) \approx \frac{N}{M} \sum_{m=1}^M E_{q(\theta)}[\log p(y_m \mid x_m, \theta)] - KL[q(\theta) \parallel p(\theta)]$$

# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:

- Variational inference:  $\phi^* = \operatorname{argmax} L(\phi)$

$$L(\phi) = E_{q_\phi(\theta)}[\log p(D \mid \theta)] - KL[q_\phi(\theta) \parallel p(\theta)]$$

- First scalable technique: **Stochastic optimization**

- i.i.d. assumption of data:  $\log p(D \mid \theta) = \sum_{n=1}^N \log p(y_n \mid x_n, \theta)$

- Enable mini-batch training with  $\{(x_m, y_m)\} \sim D^M$ :

$$L(\phi) \approx \boxed{\frac{N}{M}} \sum_{m=1}^M E_{q(\theta)}[\log p(y_m \mid x_m, \theta)] - KL[q(\theta) \parallel p(\theta)]$$

reweighting to ensure calibrated  
posterior concentration

# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:
  - 2nd scalable technique: **Monte Carlo sampling**
    - $E_{q(\theta)}[\log p(y | x, \theta)]$  intractable even with Gaussian  $q(\theta)$
    - **Solution: Monte Carlo estimate:**

$$E_{q(\theta)}[\log p(y | x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y | x, \theta_k), \quad \theta_k \sim q(\theta)$$

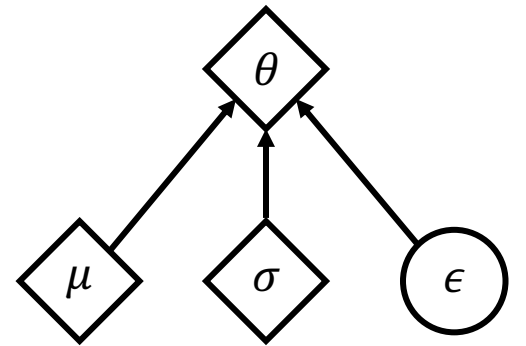


# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:
  - 2nd scalable technique: **Monte Carlo sampling**
    - $E_{q(\theta)}[\log p(y | x, \theta)]$  intractable even with Gaussian  $q(\theta)$
    - **Solution: Monte Carlo estimate:**

$$E_{q(\theta)}[\log p(y | x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y | x, \theta_k), \quad \theta_k \sim q(\theta)$$

- **Reparameterization trick** to sample mean-field Gaussians:  
 $\theta_k \sim q(\theta) \Leftrightarrow \theta_k = m_\theta + \sigma_\theta \odot \epsilon_k, \quad \epsilon_k \sim N(0, I)$



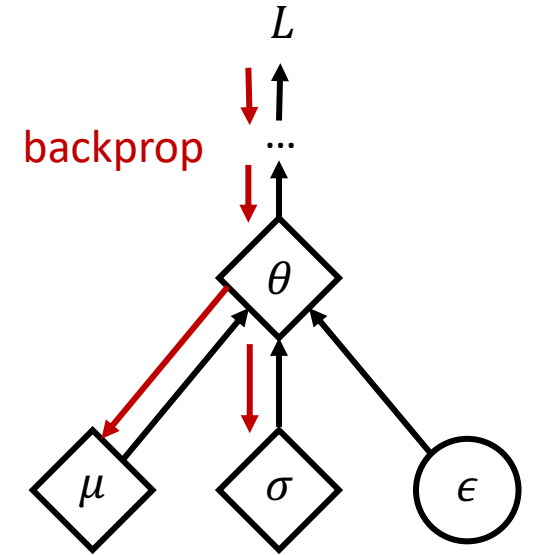
# Approximate Inference in BNNs

- Step 2: fit the  $q(\theta)$  distribution:
  - 2nd scalable technique: **Monte Carlo sampling**
    - $E_{q(\theta)}[\log p(y | x, \theta)]$  intractable even with Gaussian  $q(\theta)$
    - **Solution: Monte Carlo estimate:**

$$E_{q(\theta)}[\log p(y | x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y | x, \theta_k), \quad \theta_k \sim q(\theta)$$

- **Reparameterization trick** to sample mean-field Gaussians:  
 $\theta_k \sim q(\theta) \Leftrightarrow \theta_k = m_\theta + \sigma_\theta \odot \epsilon_k, \epsilon_k \sim N(0, I)$

$$\Rightarrow E_{q(\theta)} [\log p(y | x, \theta)] \approx \frac{1}{K} \sum_k^K \log p(y | x, \theta_k = m_\theta + \sigma_\theta \epsilon_k), \epsilon_k \sim N(0, I)$$



# Approximate Inference in BNNs

- Combining both steps:

$$L(\phi) \approx \frac{N}{M} \sum_{m=1}^M \frac{1}{K} \sum_{k=1}^K \log p(y_m | x_m, \theta_k) - \underline{KL[q(\theta) || p(\theta)]}, \theta_k \sim q(\theta)$$

analytic between two Gaussians

(if not, can also be estimated with Monte Carlo)

In regression:

$$p(y | x, \theta) = N(f_{\theta}(x), \sigma^2)$$

In classification:

$$p(y | x, \theta) = \text{Categorical}(\text{logit} = f_{\theta}(x))$$

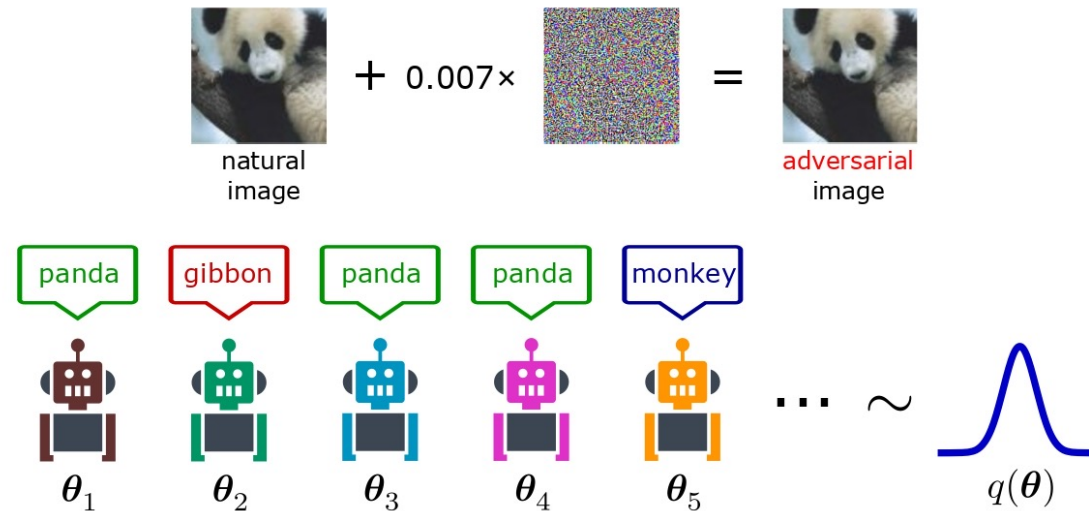
# Approximate Inference in BNNs

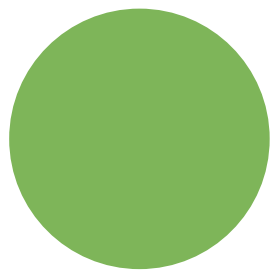
- Step 3: compute prediction with Monte Carlo approximations:

$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \underline{\theta_k \sim q(\theta)}$$

Mean-field Gaussian case:

$$\theta_k = m_\theta + \sigma_\theta \odot \epsilon_k, \quad \epsilon_k \sim N(0, I)$$





# Part II: Bayesian MLPs

- **Implement various BNN methods for MLP architectures**
- **Regression example test**
- **Case study 1: Bayesian Optimisation with UCB**



<https://bit.ly/3zF1zvA>

Instructions for using this Google Colab notebook:

- Make sure you have signed in with your Google account;
- Click “File > Save a copy in Drive” to create your own copy;
- Let’s play around with the demo using your own copy!

# Findings with MFVI for Bayesian MLPs

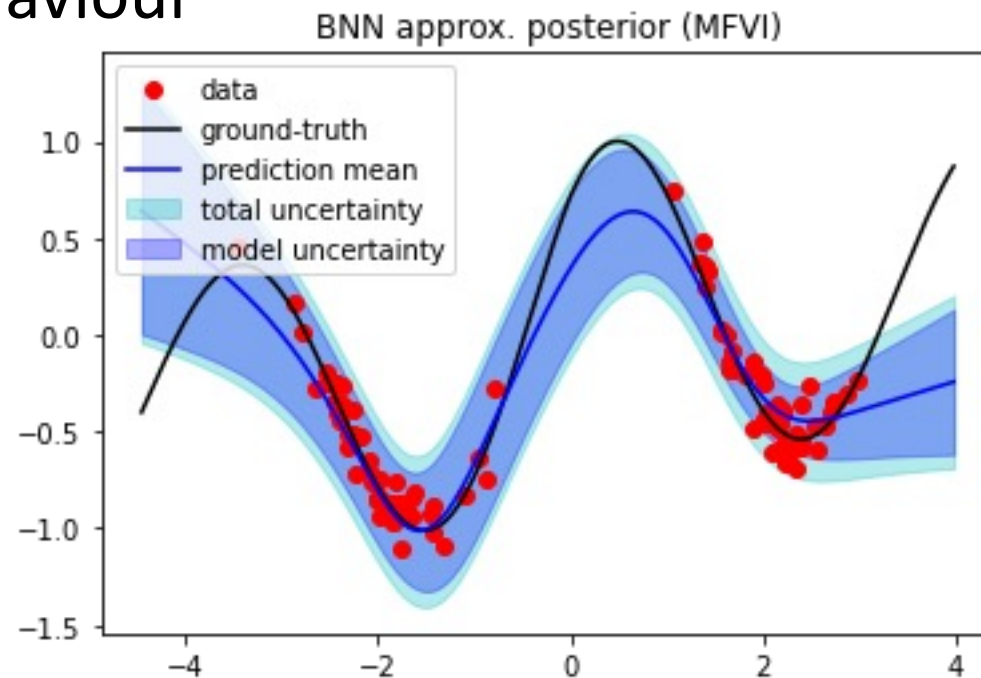


# Findings with MFVI for Bayesian MLPs

- MFVI tends to underfit
  - Initialisation matters
  - Tuning the beta parameter also helps

# Findings with MFVI for Bayesian MLPs

- MFVI tends to underfit
  - Initialisation matters
  - Tuning the beta parameter also helps
- Uncertainty behaviour



# Using other $q$ distributions?

- Using more complicated  $q$  distributions?
  - Pros: more flexible approximations  $\Rightarrow$  better posterior approximations (?)
  - Cons: higher time & space complexities

# Using other $q$ distributions?

- Using more complicated  $q$  distributions?
  - Pros: more flexible approximations  $\Rightarrow$  better posterior approximations (?)
  - Cons: higher time & space complexities
- We will look at 2 alternatives:
  - “Last-layer BNN”: Full covariance Gaussian approximations for the last layer
  - MC-Dropout: adding dropout layers and run them in both train & test time

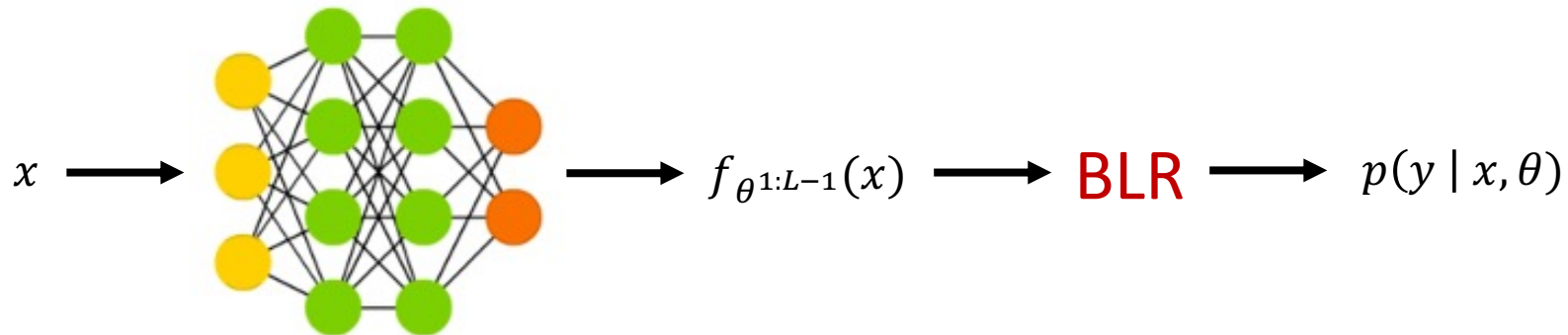
# “Last-layer BNN”

- Use deterministic layers for all but the last layer
- For the last layer: Use Full-covariance Gaussian approximate posterior:

$$q(\theta^l) = \delta(W^l = M^l, b^l = m^l), l = 1, \dots, L - 1,$$

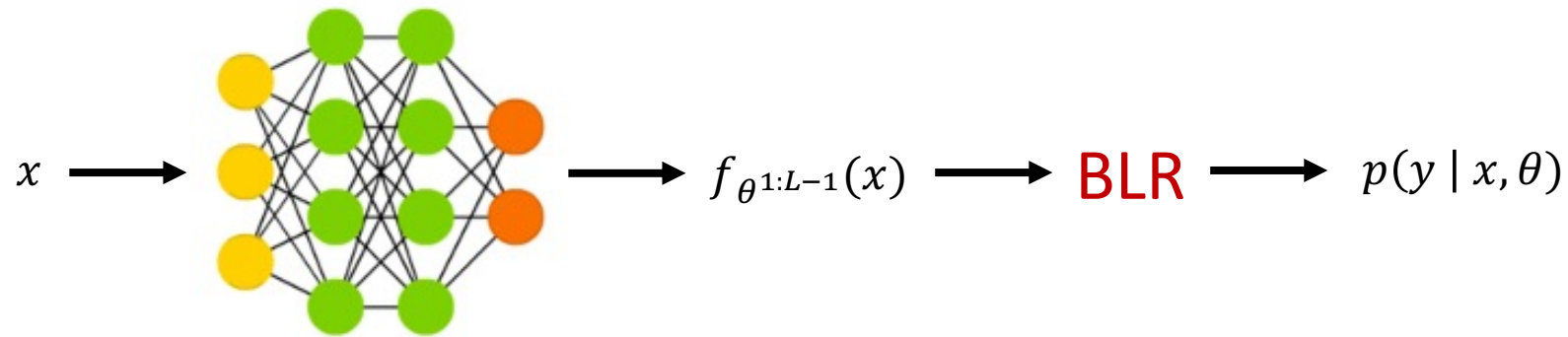
$$q(\theta^L) = N(\text{vec}(\theta^L); \text{vec}(\mu^L), \Sigma), \theta^L = \{W^L, b^L\}$$

- For regression this is equivalent to **Bayesian linear regression (BLR)** with **NN-based non-linear features**



# “Last-layer BNN”

- Use deterministic layers for all but the last layer
- For the last layer: Use Full-covariance Gaussian approximate posterior
- For regression this is equivalent to **Bayesian linear regression (BLR)** with **NN-based non-linear features**



$$L = E_q[\log p(D | \theta)] - KL[q(\theta^L) \parallel p(\theta^L)]$$

Use deterministic weights for  $l = 1, \dots, L - 1$   
Sample  $W^L \sim q(W^L)$

KL regulariser for the last layer only

# MC-Dropout

- Add dropout layers to the network
- Perform dropout during training

$$L = E_q[\log p(D|\theta)] - (1 - \pi)\ell_2(\phi)$$

The MC sampling procedure is implicitly defined

L2 regulariser on the variational parameters

Dropout rate

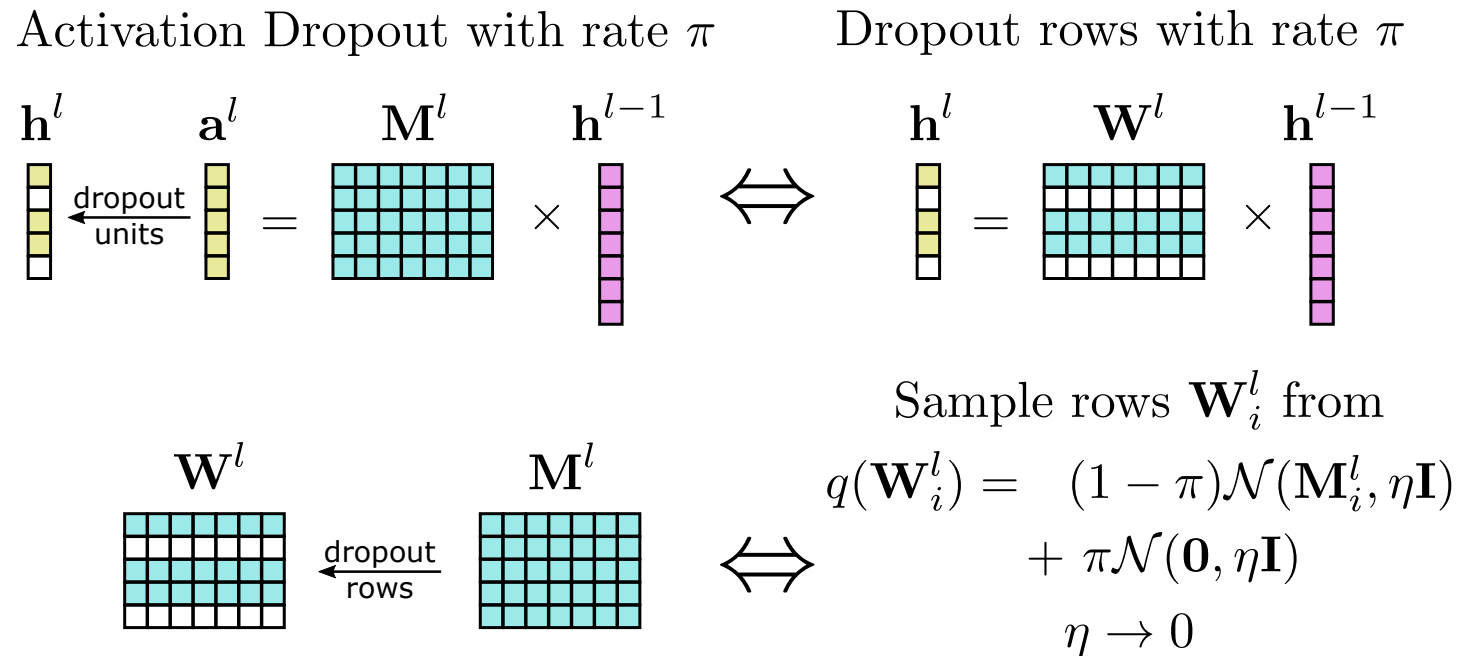
- In test time, run multiple forward passes with dropout

$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$

The MC sampling procedure is implicitly defined

# MC-Dropout

- Two equivalent ways to implement MC-Dropout:



(Similar logic applies when including the bias terms, see lecture notes.)

(Notice that pytorch's nn.Linear layer uses formats like  $xW^T$  instead of  $Wx$ .)



# Using other $q$ distributions?

- What you'll do for the next part of the tutorial:
  - Implement MC-Dropout in 2 ways
  - Run the regression sample with the 2 approximation methods discussed
  - Compare with MFVI

# Case study 1: Bayesian Optimisation

- Imagine you'd like to solve the following task:

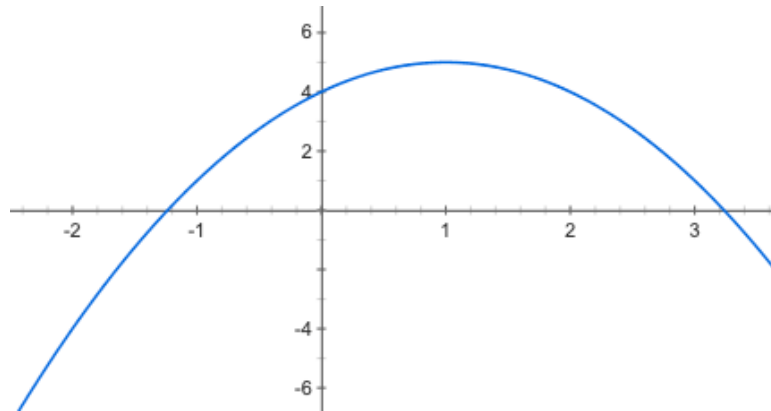
$$x^* = \operatorname{argmax}_x f_0(x)$$

# Case study 1: Bayesian Optimisation

- Imagine you'd like to solve the following task:

$$x^* = \operatorname{argmax}_x f_0(x)$$

Known functional form of  $f_0$ :



Gradient descent, Newton's method,

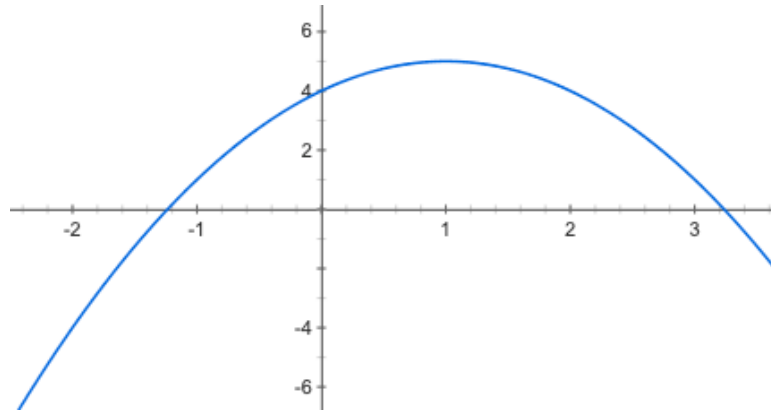
...

# Case study 1: Bayesian Optimisation

- Imagine you'd like to solve the following task:

$$x^* = \operatorname{argmax}_x f_0(x)$$

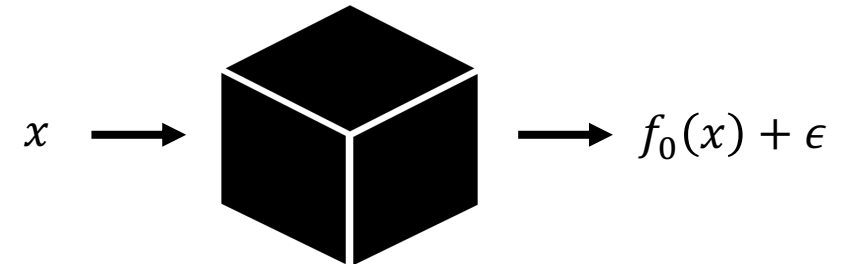
Known functional form of  $f_0$ :



Gradient descent, Newton's method,

...

**Unknown** functional form of  $f_0$ :

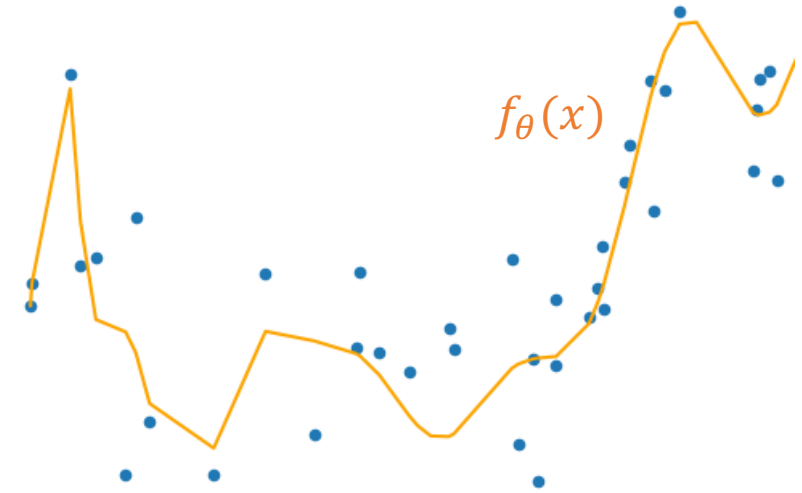
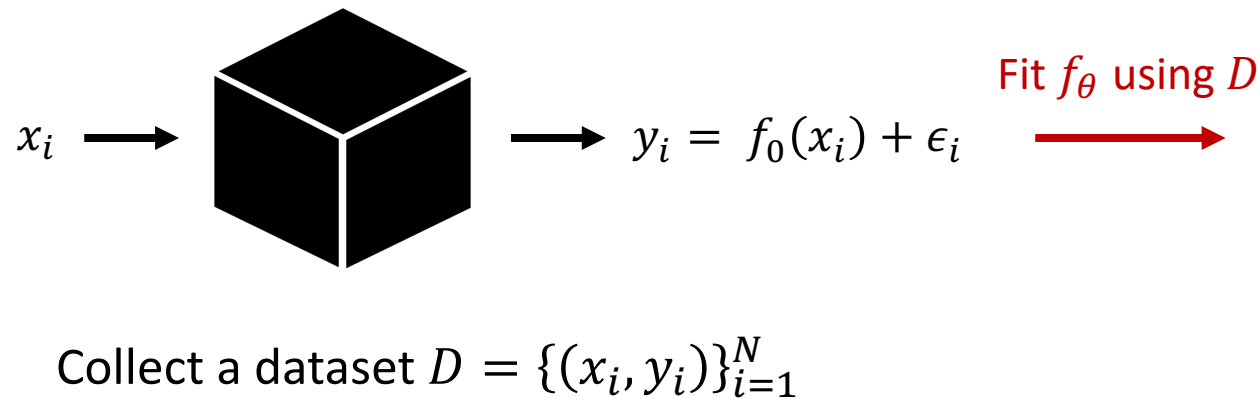


(can only query (noisy) function values)



# Case study 1: Bayesian Optimisation

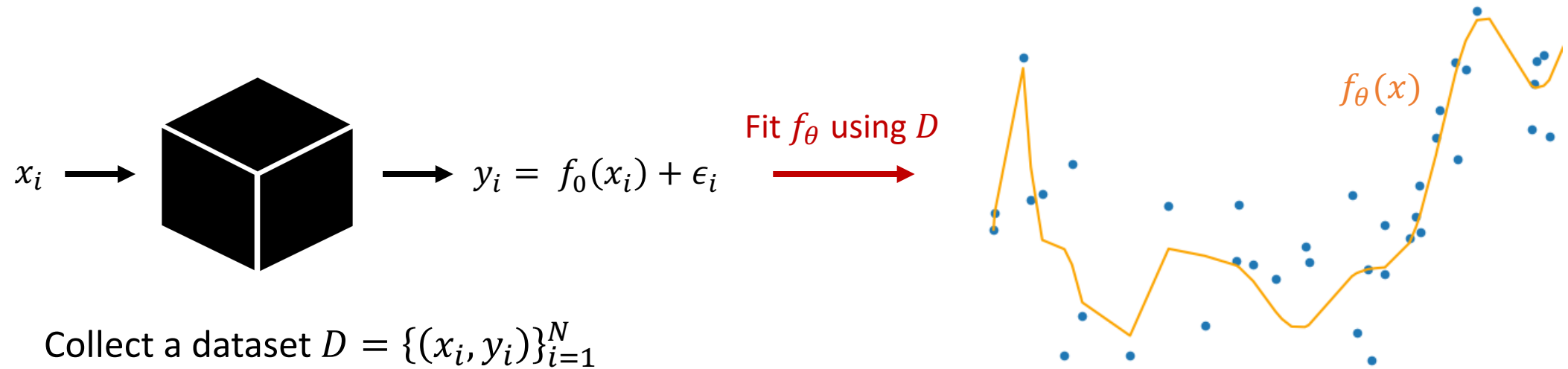
- Idea 1: fit a surrogate function  $f_\theta \approx f_0$



$f_\theta(x)$  has a known (parametric) form  
 $\Rightarrow$  find maximum using e.g., Newton's method

# Case study 1: Bayesian Optimisation

- Idea 1: fit a surrogate function  $f_\theta \approx f_0$

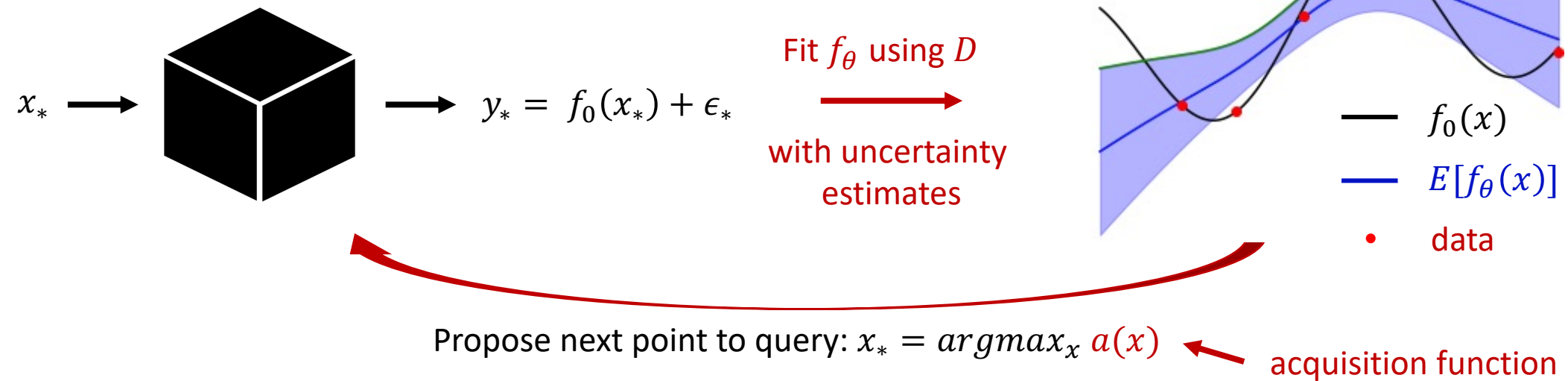


- Issues of this approach:
  - Need to collect a lot of datapoints for accurate fitting of  $f_\theta$
  - Do not consider uncertainty at unseen locations

# Case study 1: Bayesian Optimisation

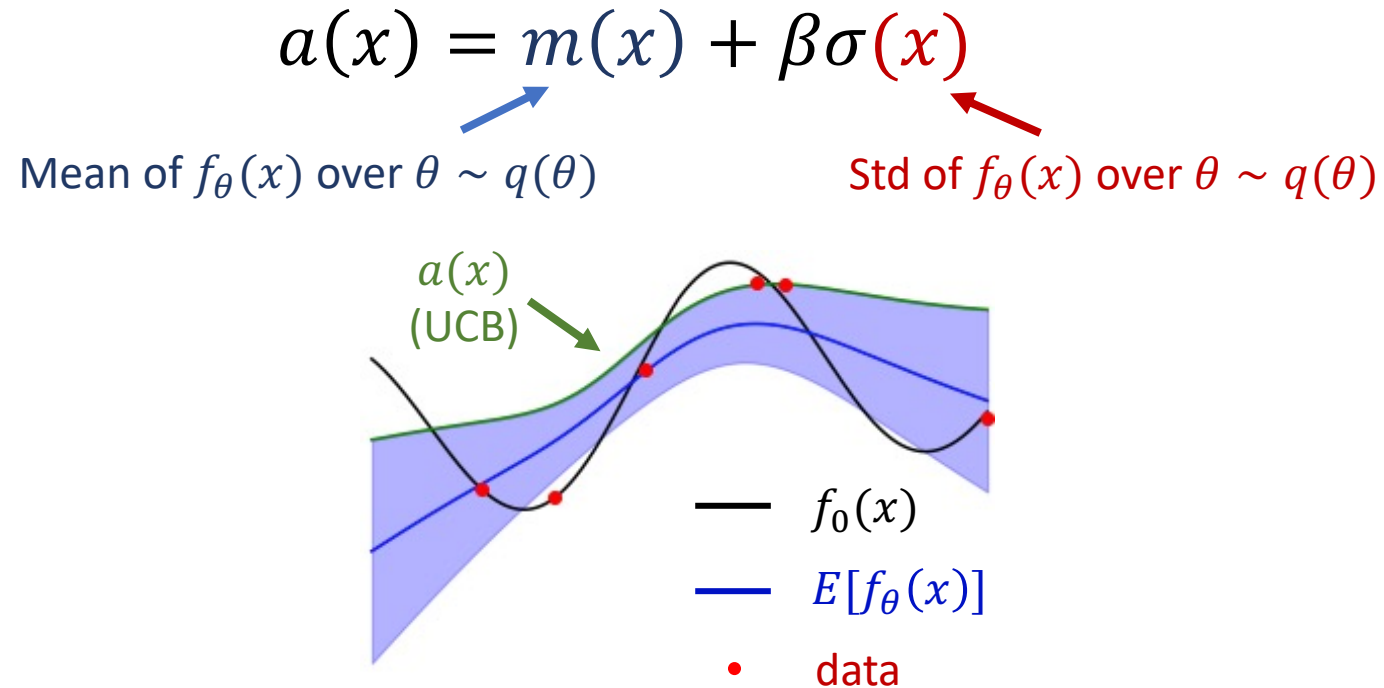
- Idea of BO: iterate the following steps
  - fit a surrogate function  $f_\theta$  **with uncertainty estimates**
  - Use the surrogate function to **guide the dataset collection process**

Update the dataset  $D = D \cup \{(x_*, y_*)\}$



# Case study 1: Bayesian Optimisation

- Upper confidence bound (UCB): a widely used acquisition function



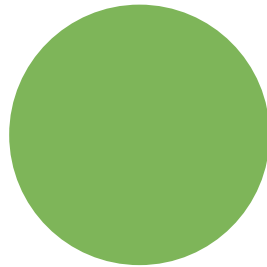


# Case study 1: Bayesian Optimisation

- What you'll do for the case study part of the tutorial:
  - Implement UCB acquisition function
  - Run the BO example
  - Play around with hyper-parameters and other settings

Example answers of the tutorial demos:

Regression: <https://bit.ly/39eZHit>



# Part III: Bayesian ConvNets

- **Classification example test**
- **Case study 2: Detecting adversarial examples**



<https://bit.ly/3Hd1Ass>

Instructions for using this Google Colab notebook:

- Make sure you have signed in with your Google account;
- Click “File > Save a copy in Drive” to create your own copy;
- Use GPU: in “Runtime > Change runtime type”, choose “GPU” for “hardware accelerator”
- Let’s play around with the demo using your own copy!

# Case study 2: Detecting adversarial examples

- Hypothesis:
  - Adversarial examples are regarded as OOD data
  - BNNs become uncertain about their prediction on OOD data
  - $\Rightarrow$  uncertainty measures can be used for detecting adversarial examples

# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

Due to lack of “knowledge”  
(reducible when having more data)

Due to inherent stochasticity in data  
(non-reducible)

Imagine flipping a coin:

- **Epistemic uncertainty**: “How much do I believe the coin is fair?”
  - Model’s belief after seeing the population
  - Reduces when having more data
- **Aleatoric uncertainty**: “What’s the next coin flip outcome?”
  - Individual experiment outcome
  - Non-reducible



# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

Due to lack of “knowledge”  
(reducible when having more data)

Due to inherent stochasticity in data  
(non-reducible)

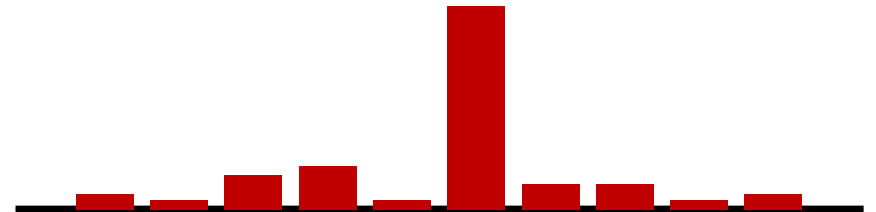
Computing uncertainty in classification models:

$$H[p] = -\sum_{c=1}^C p_c \log p_c, \quad p = (p_1, \dots, p_C), \sum_{c=1}^C p_c = 1$$

High entropy:  $H[p] \rightarrow \log C$





Low entropy:  $H[p] \rightarrow 0$



# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

  
Due to lack of “knowledge”  
(reducible when having more data)

  
Due to inherent stochasticity in data  
(non-reducible)

Computing uncertainty in classification models:

Recall for Bayesian predictive distribution:

$$p(y^* | x^*, D) = \int p(y^* | x^*, \theta) p(\theta | D) d\theta$$

$$H[y^* | x^*, D] = I[y^*; \theta | x^*, D] + E_{p(\theta | D)}[H[y^* | x^*, \theta]]$$

Total entropy of the  
predictive distribution


Mutual information  
between  $y^*$  and  $\theta$


Conditional entropy  
under posterior



# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

  
Due to lack of “knowledge”  
(reducible when having more data)

  
Due to inherent stochasticity in data  
(non-reducible)

Computing uncertainty in classification models:

Recall for Bayesian predictive distribution **with approximation**:

$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$


Total entropy (for total uncertainty):

$$H[y^* | x^*, D] \approx H\left[\frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k)\right]$$

# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

  
Due to lack of “knowledge”  
(reducible when having more data)

  
Due to inherent stochasticity in data  
(non-reducible)

Computing uncertainty in classification models:

Recall for Bayesian predictive distribution **with approximation**:


$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$


Conditional entropy (for aleatoric uncertainty):

$$E_{p(\theta | D)}[H[y^* | x^*, \theta]] \approx \frac{1}{K} \sum_{k=1}^K H[p(y^* | x^*, \theta_k)]$$

# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

  
Due to lack of “knowledge”  
(reducible when having more data)

  
Due to inherent stochasticity in data  
(non-reducible)

Computing uncertainty in classification models:

Recall for Bayesian predictive distribution **with approximation**:


$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$


**Mutual information (for epistemic uncertainty):**

$$I[y^*; \theta | x^*, D] \approx H\left[\frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k)\right] - \frac{1}{K} \sum_{k=1}^K H[p(y^* | x^*, \theta_k)]$$

# Uncertainty measures

Total uncertainty = **epistemic uncertainty** + **aleatoric uncertainty**

  
Due to lack of “knowledge”  
(reducible when having more data)

  
Due to inherent stochasticity in data  
(non-reducible)

Computing uncertainty in classification models:

Recall for Bayesian predictive distribution **with approximation**:

$$p(y^* | x^*, D) \approx \frac{1}{K} \sum_{k=1}^K p(y^* | x^*, \theta_k), \quad \theta_k \sim q(\theta)$$

**Mutual information (for epistemic uncertainty) if you can do exact inference:**

$$I[y^*; \theta | x^*, D] = E_{p(y^* | x^*, D)} [KL[p(\theta | D, x^*, y^*) \parallel p(\theta | D)]]$$

“What the model thinks the posterior is going to change if we add new observation at location  $x^*$ ”

# Case study 2: Detecting adversarial examples

- What you'll do for the case study part of the tutorial:
  - Implement the uncertainty measures
    - Total entropy, conditional entropy, and mutual info
  - Run adversarial attacks on various trained networks
  - See how diversity helps in detecting adversarial examples
    - Detection by thresholding the uncertainty measures
    - We consider best TPR with  $\text{FPR} \leq 5\%$

# Ensemble BNNs

- Define  $q$  distribution as mixture of mean-field Gaussian:

$$q(\theta) = \frac{1}{S} \sum_{s=1}^S q(\theta|s), \quad q(\theta|s) = N(\theta; \mu_s, \text{diag}(\sigma_s^2))$$

- Objective is still a valid lower-bound to  $\log p(D)$ :

$$L = \frac{1}{S} \sum_{s=1}^S ELBO[q(\theta|s)], \quad ELBO[q(\theta|s)] = E_{q(\theta|s)}[\log p(D | \theta)] - KL[q(\theta|s) \parallel p(\theta)]$$

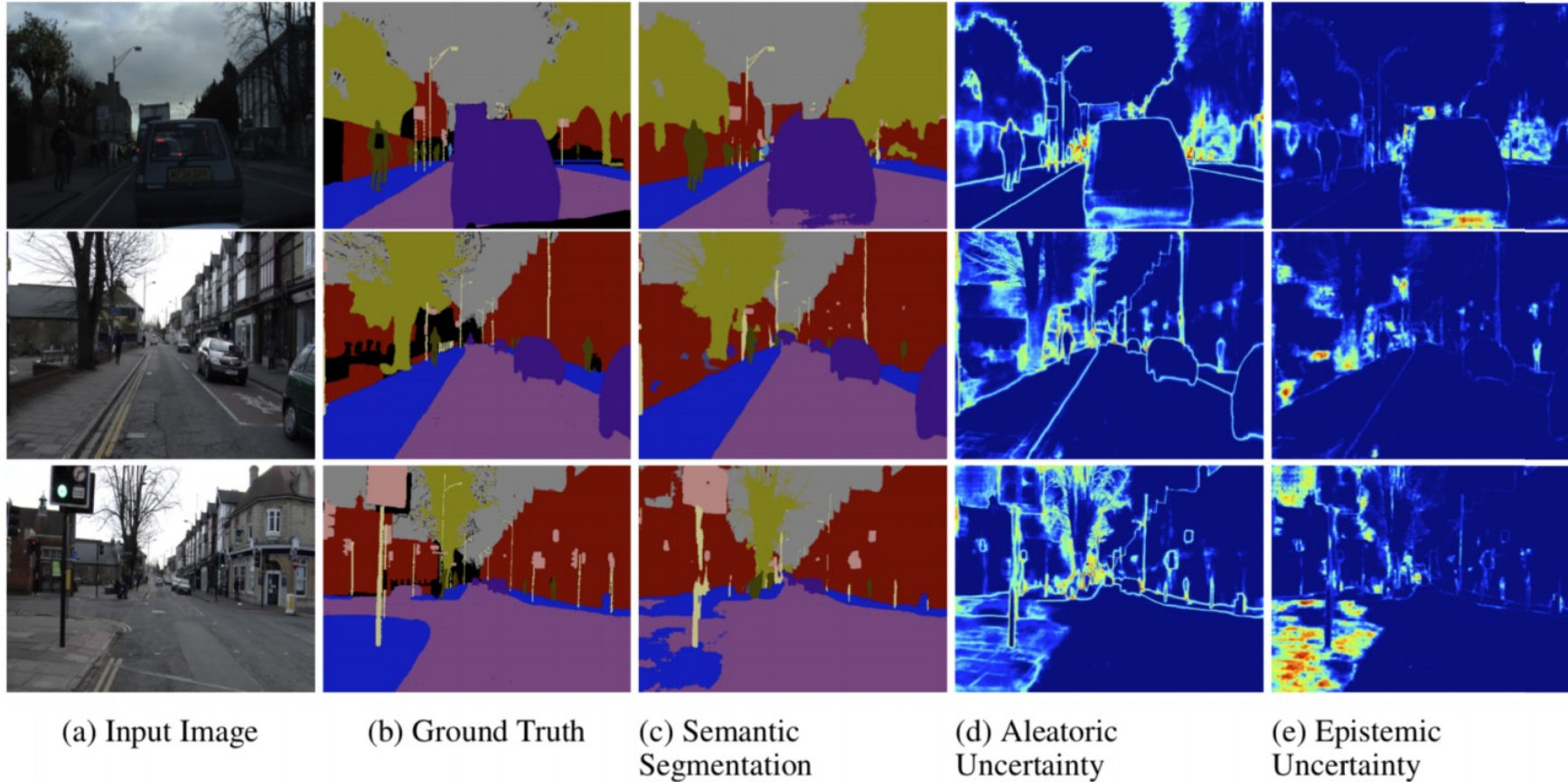
- The parameters of  $q(\theta|s)$  for different  $s$  are **independent**  
 $\Rightarrow$  train  $S$  number of MFVI-BNNs independently



# Part IV: Advances & Future Works

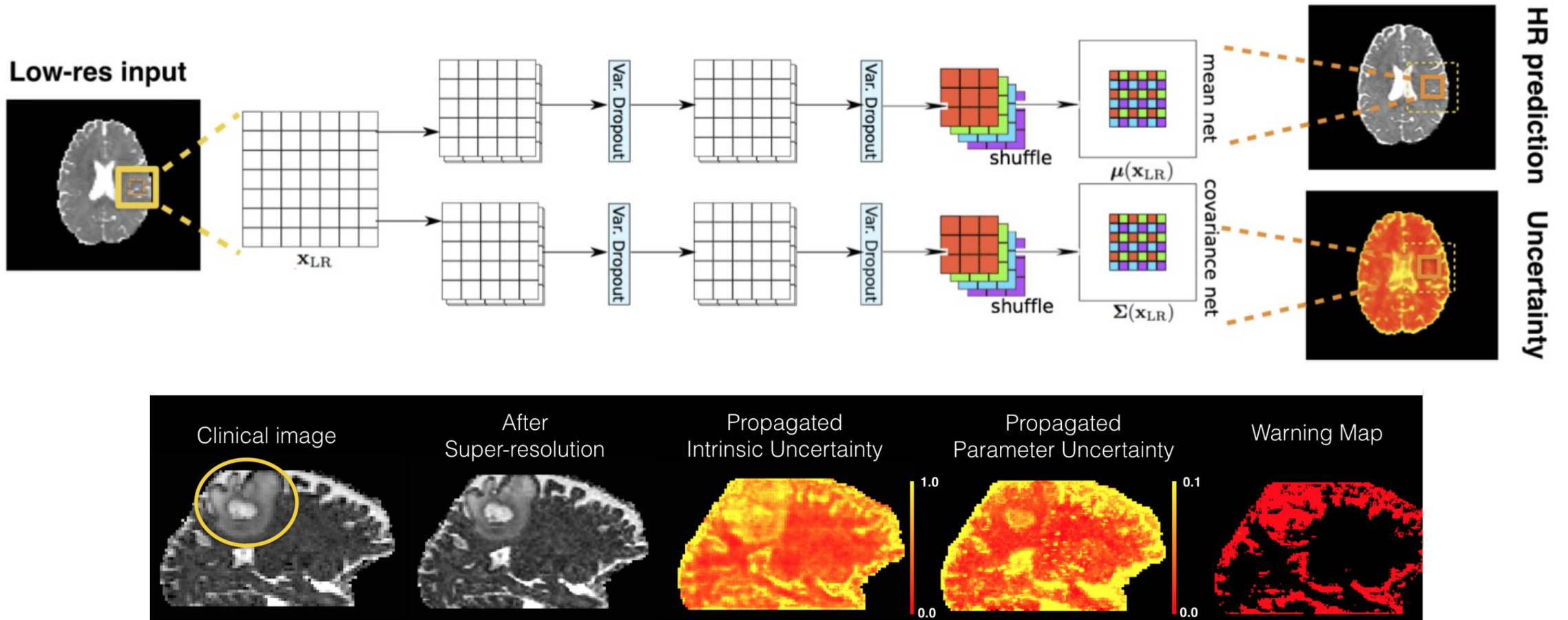
- **Various applications**
- **Overview of recent progresses**
- **Future directions**

# Applications of BNNs: Image Segmentation

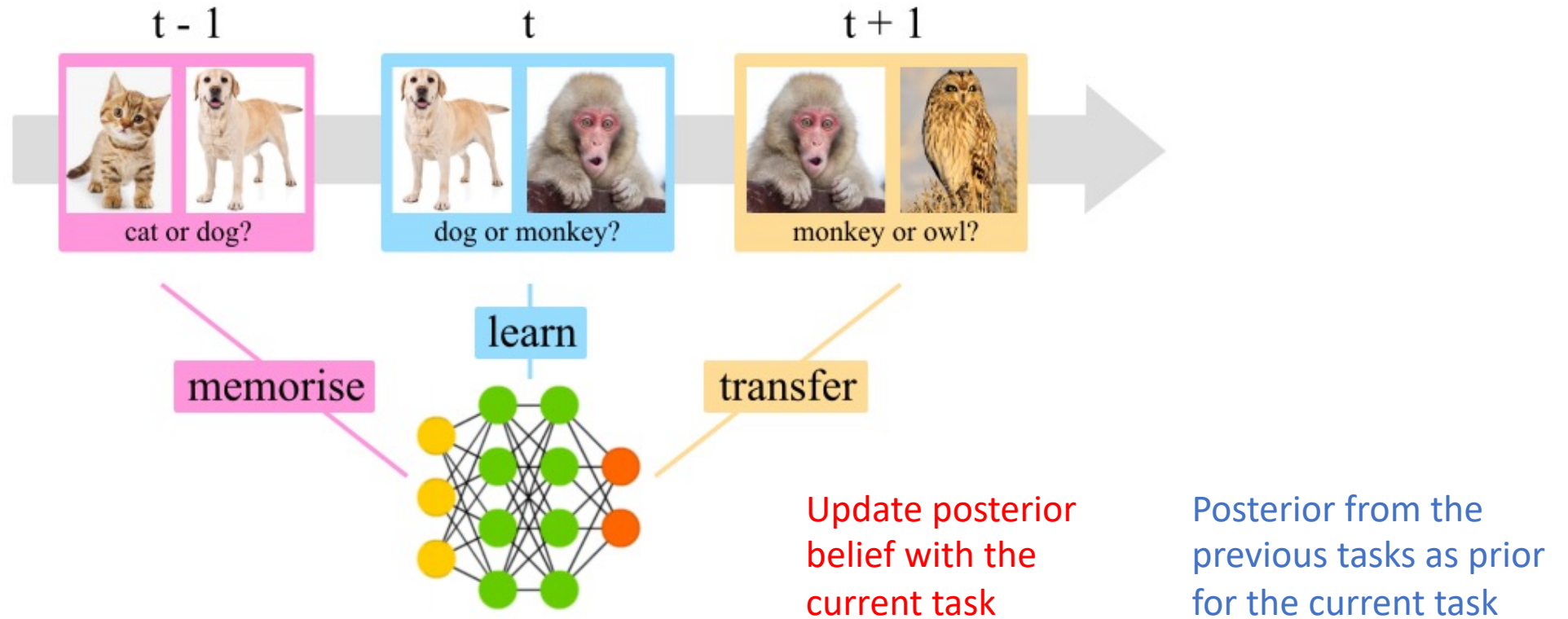




# Applications of BNNs: Super Resolution



# Applications of BNNs: Continual Learning



$$L_{VCL}^t(q_t(\theta)) = E_{q_t(\theta)}[\log p(D_t | \theta)] - KL[q_t(\theta) || q_{t-1}(\theta)]$$

# Recent Progress in BNNs: Inference

$$\text{SGD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$$

$$\text{SGLD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim N(0, I)$$

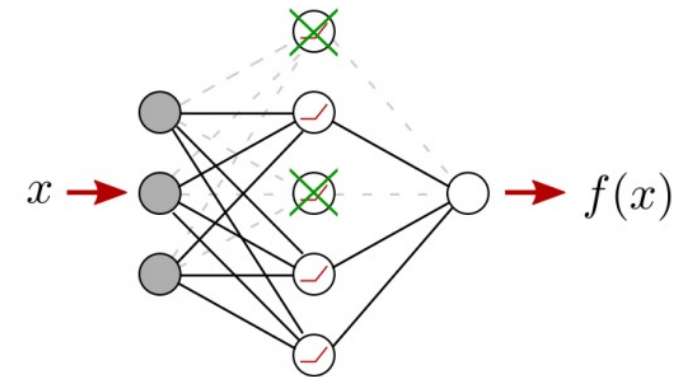
Stochastic gradient MCMC

# Recent Progress in BNNs: Inference

$$\text{SGD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$$

$$\text{SGLD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim N(0, I)$$

Stochastic gradient MCMC



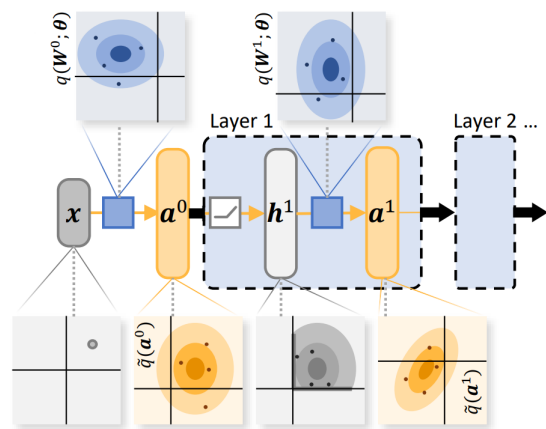
Monte Carlo dropout

# Recent Progress in BNNs: Inference

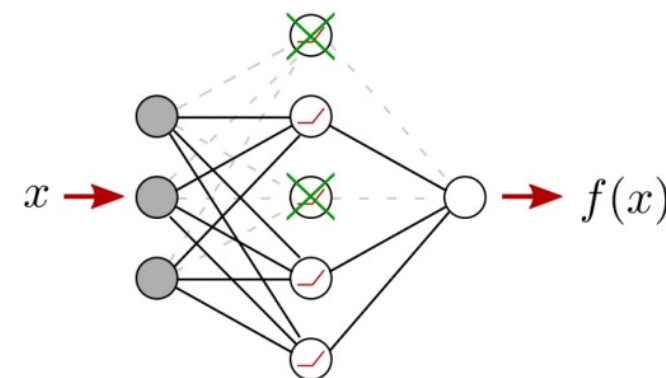
$$\text{SGD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$$

$$\text{SGLD: } \theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim N(0, I)$$

Stochastic gradient MCMC



Deterministic approximations



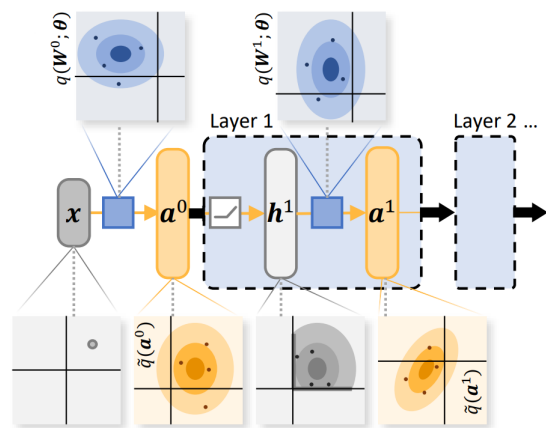
Monte Carlo dropout

# Recent Progress in BNNs: Inference

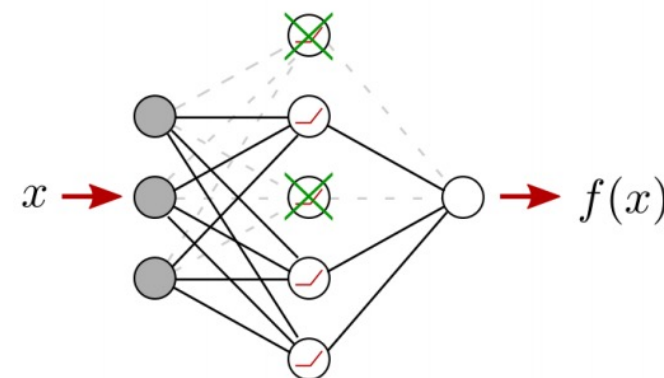
SGD:  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t)$

SGLD:  $\theta_{t+1} = \theta_t - \eta \nabla_{\theta_t} \tilde{U}(\theta_t) + \sqrt{2\eta} \epsilon, \quad \epsilon \sim N(0, I)$

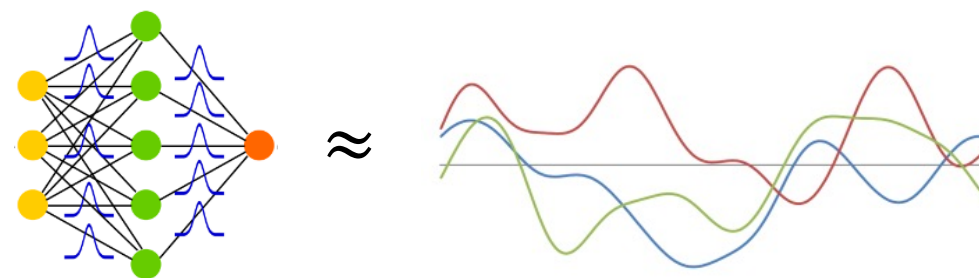
Stochastic gradient MCMC



Deterministic approximations

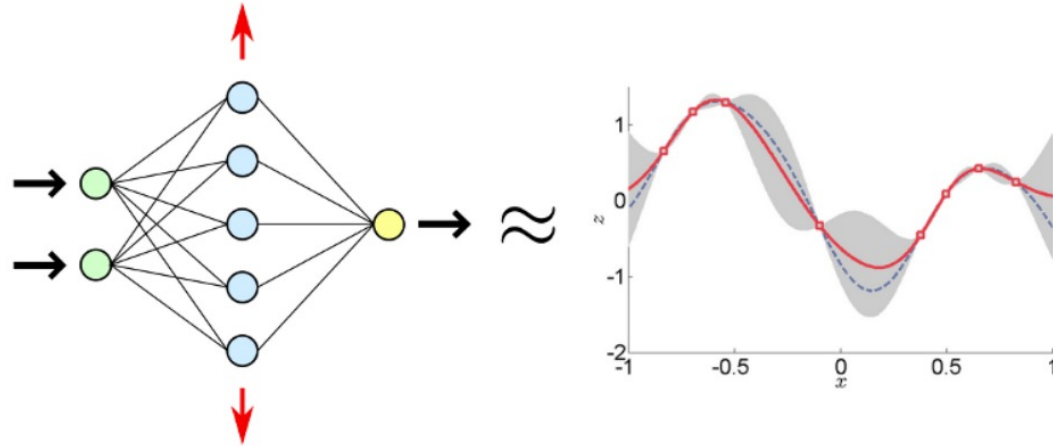


Monte Carlo dropout



Function space approximate inference

# Recent Progress in BNNs: Theory



## Connections to GPs:

- BNN with very wide hidden layers  
     $\approx$  Gaussian process
- Width limit convergence: in both  
    prior (Neal's result) and posterior

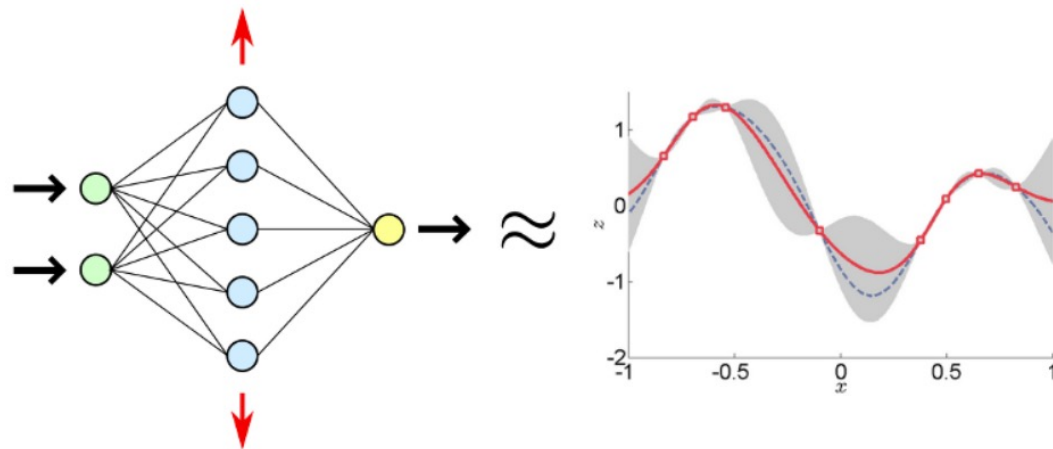
Neal. Bayesian Learning for Neural Networks. PhD Thesis, 1996

Matthews et al. Gaussian Process Behaviour in Wide Deep Neural Networks. ICLR 2018

Lee et al. Deep Neural Networks as Gaussian Processes. ICLR 2018

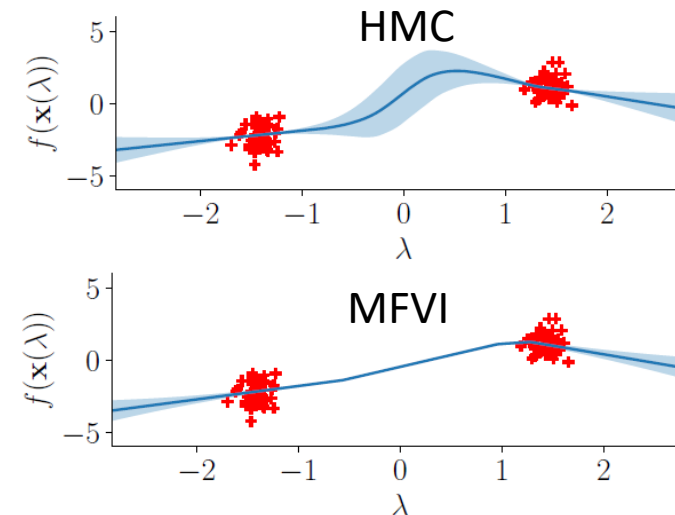
Hron et al. Exact posterior distributions of wide Bayesian neural networks. 2020

# Recent Progress in BNNs: Theory



## Connections to GPs:

- BNN with very wide hidden layers  $\approx$  Gaussian process
- Width limit convergence: in both prior (Neal's result) and posterior



## Approx. vs exact inference:

- Theoretical limitation of MFVI in shallow BNNs with ReLU activations
- Empirically deep BNNs with MFVI still fails in certain cases

Foong et al. On the Expressiveness of Approximate Inference in Bayesian Neural Networks. NeurIPS 2020

Farquhar et al. Liberty or Depth: Deep Bayesian Neural Nets Do Not Need Complex Weight Posterior Approximations. NeurIPS 2020

Coker et al. Wide Mean-Field Bayesian Neural Networks Ignore the Data. AISTATS 2022



# Future directions

- Understanding BNN behaviour:
  - How would  $q(f)$  behave given a particular form of  $q(W)$ ?
  - Is weight-space objective appropriate for MFVI?
  - We don't understand very well the optimisation properties of VI-BNN
- Computational complexity overhead: worth it?
  - How can we make the approximate posterior more efficient in both time and space complexities?
- Priors for BNNs
  - “Default” Gaussian prior  $N(W; 0, \sigma^2 I)$ : the right prior?
  - How to think about priors in function space?
- Applications
  - Improve for applications that require good uncertainty estimates



# Thank You!

Questions? Ask NOW or email:  
[yingzhen.li@imperial.ac.uk](mailto:yingzhen.li@imperial.ac.uk)

Example answers of the tutorial demos:

Regression: <https://bit.ly/39eZHit>

Classification: <https://bit.ly/3QikcLO>