

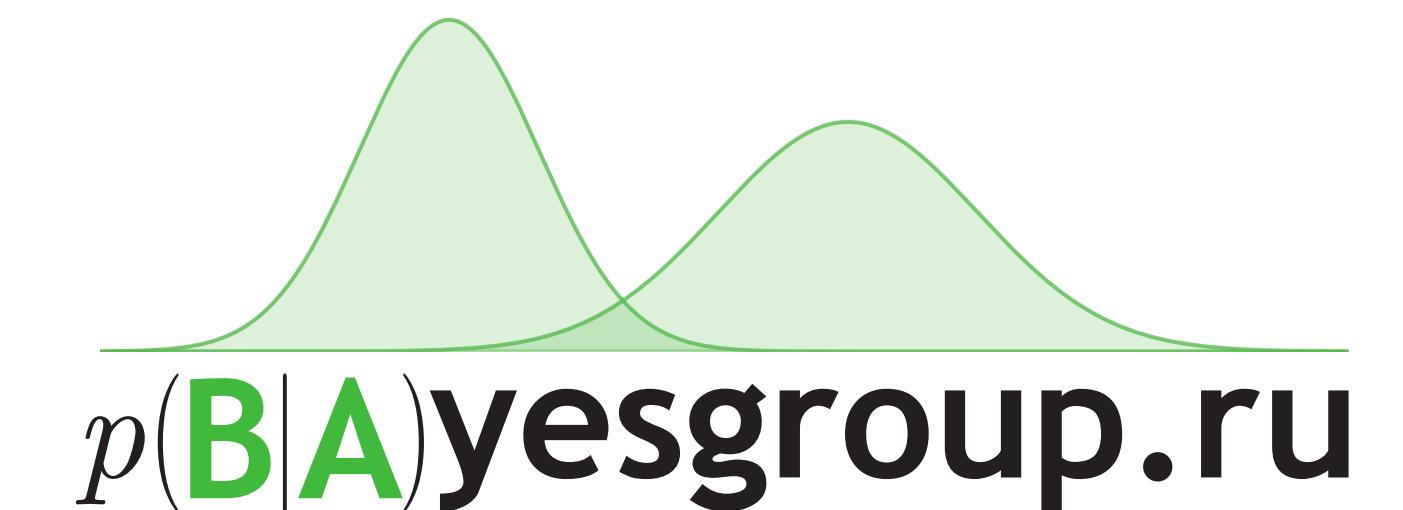
Bayesian neural networks

Nadia Chirkova

Higher School of Economics, Samsung-HSE Laboratory
Moscow, Russia



SAMSUNG
Research



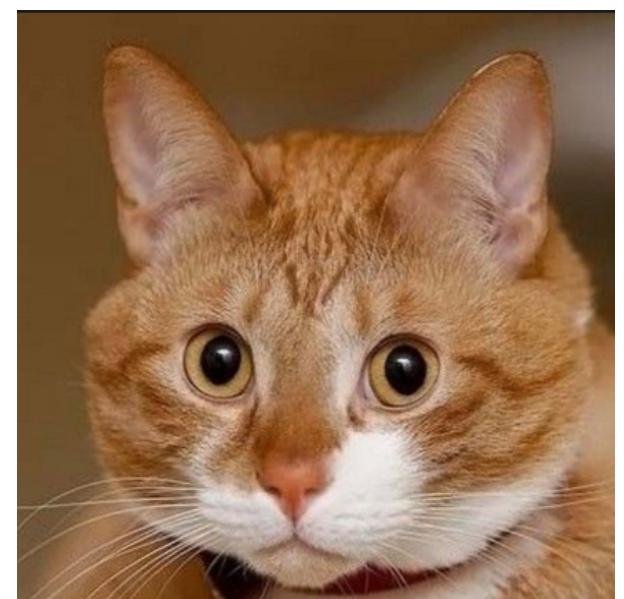
Plan

- Advantages of using Bayesian neural networks (30 min)
- Training Bayesian neural networks (30 min)
- Q&A + exercises (30 min)

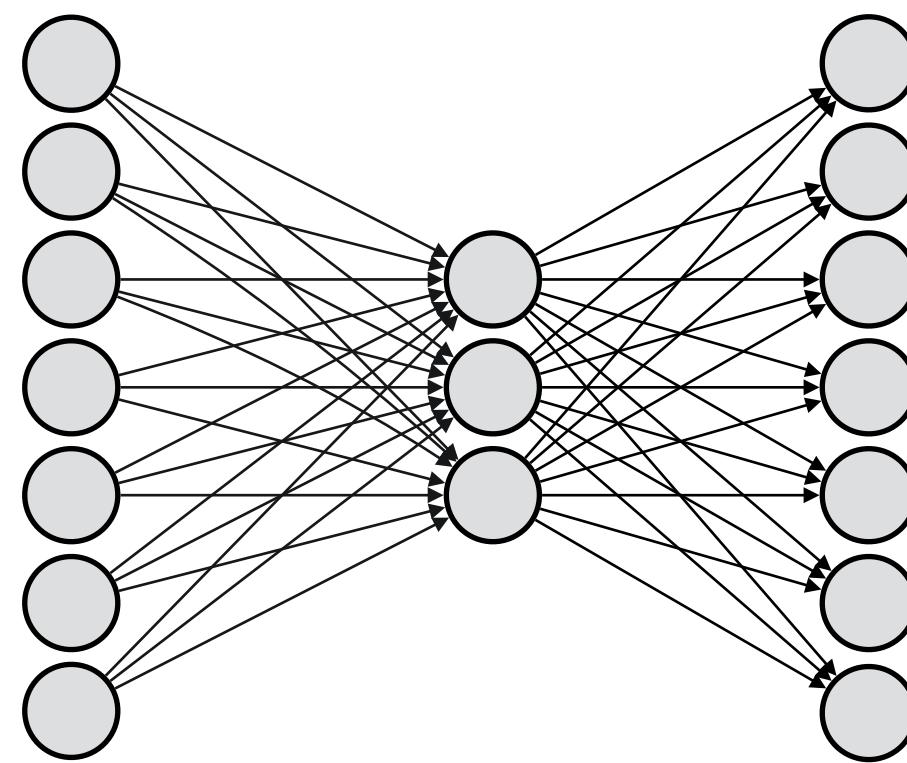
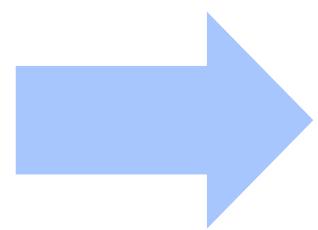
Plan

- Advantages of using Bayesian neural networks
- Training Bayesian neural networks
- Q&A + exercises

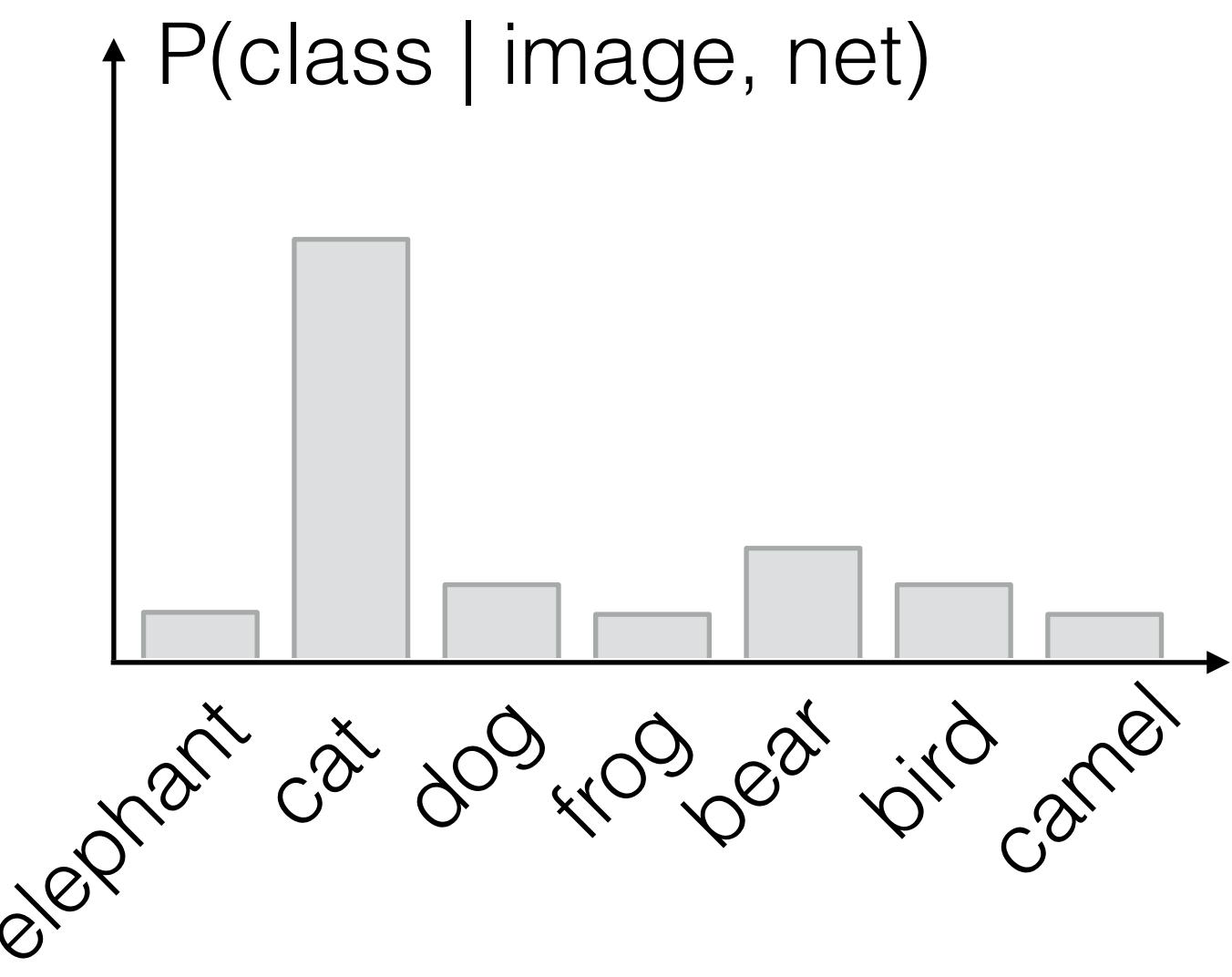
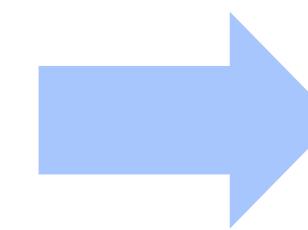
Neural networks



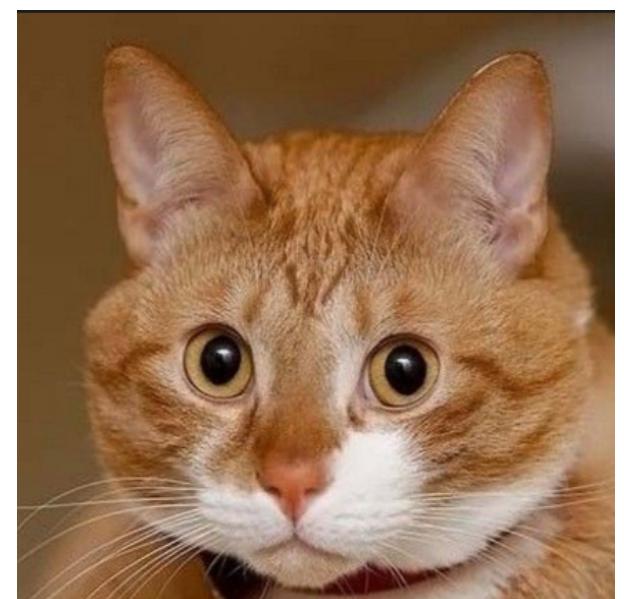
input x



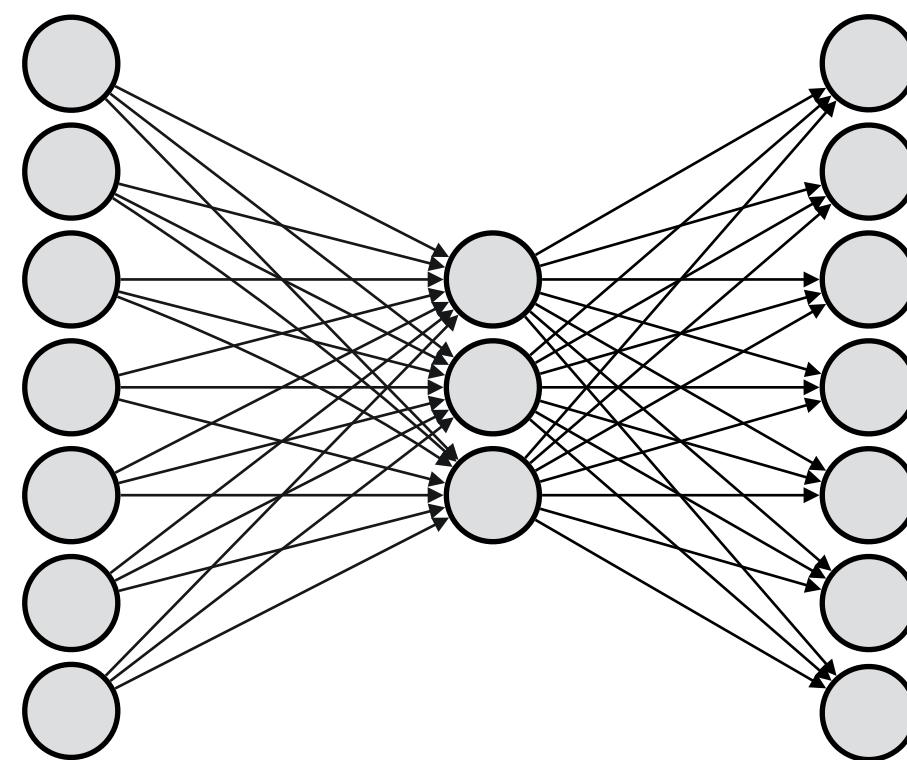
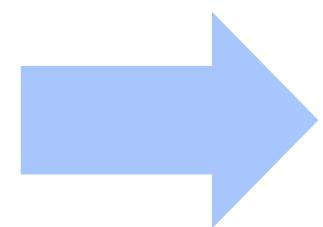
neural network
with weights w



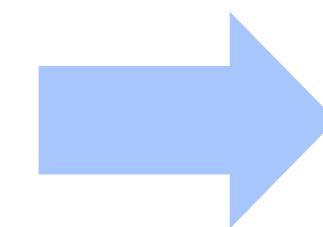
Neural networks



input x



neural network
with weights w



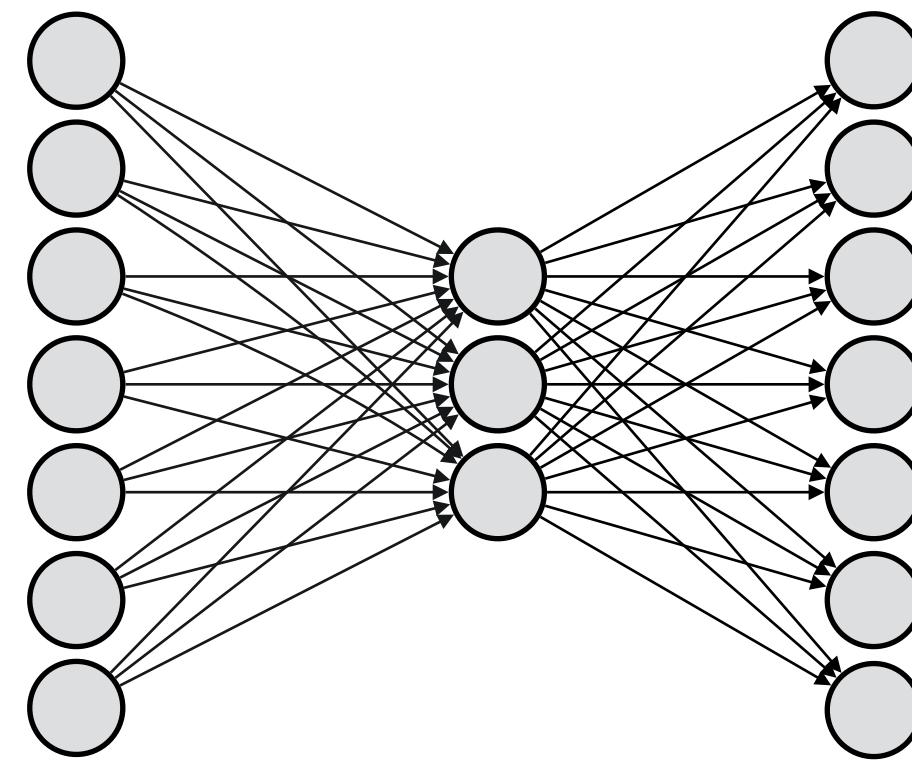
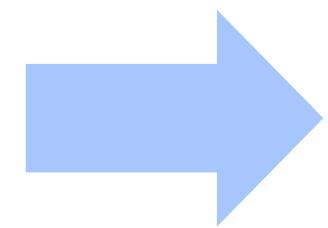
$$p(y = \text{"cat"} | x, w)$$



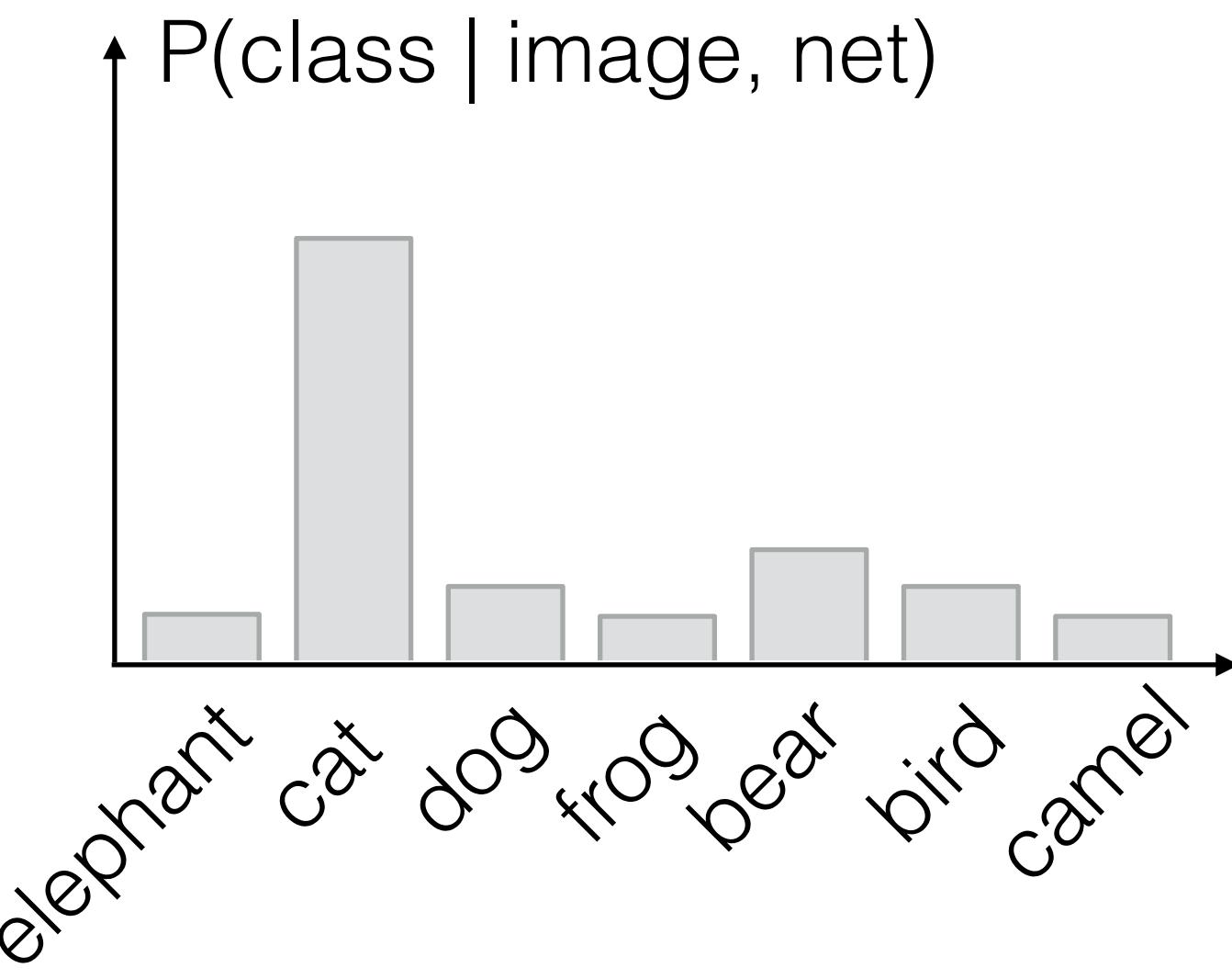
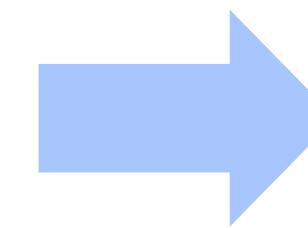
Neural networks



input x



neural network
with weights w



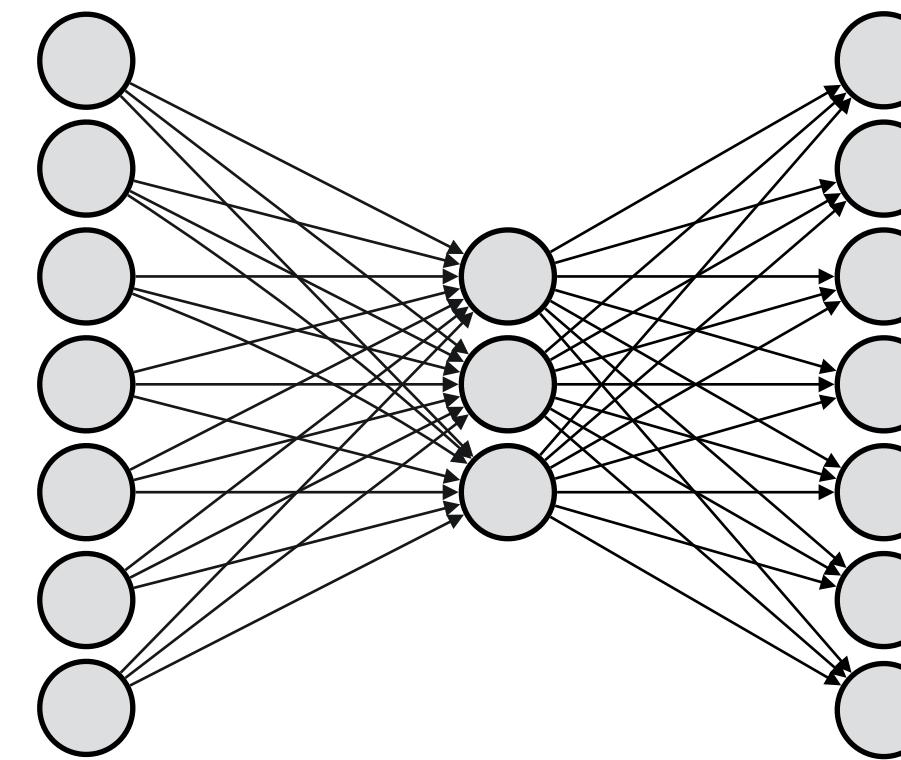
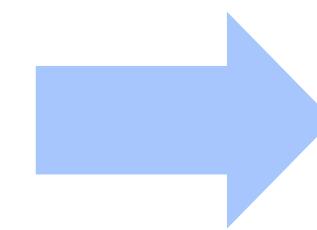
Training — optimization
over weights w
using stochastic
gradient descend:

$$-\text{DataLoss}(X, Y, w) \rightarrow \max_w$$

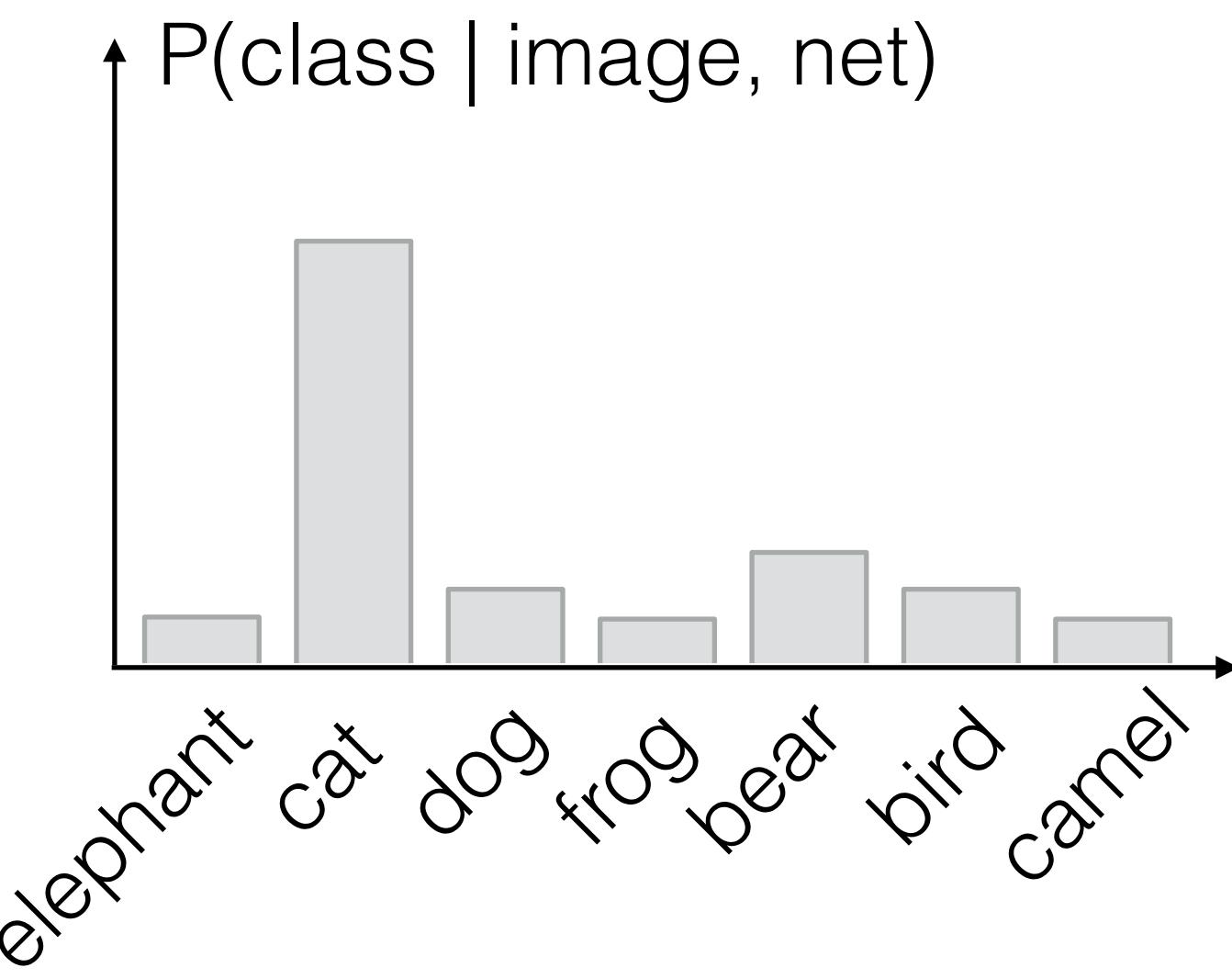
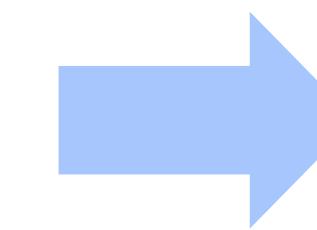
Neural networks



input x



neural network
with weights w



Training — optimization
over weights w
using stochastic
gradient descend:

$$\sum_{i=1}^N \log p(y^i|x^i, w) \rightarrow \max_w$$

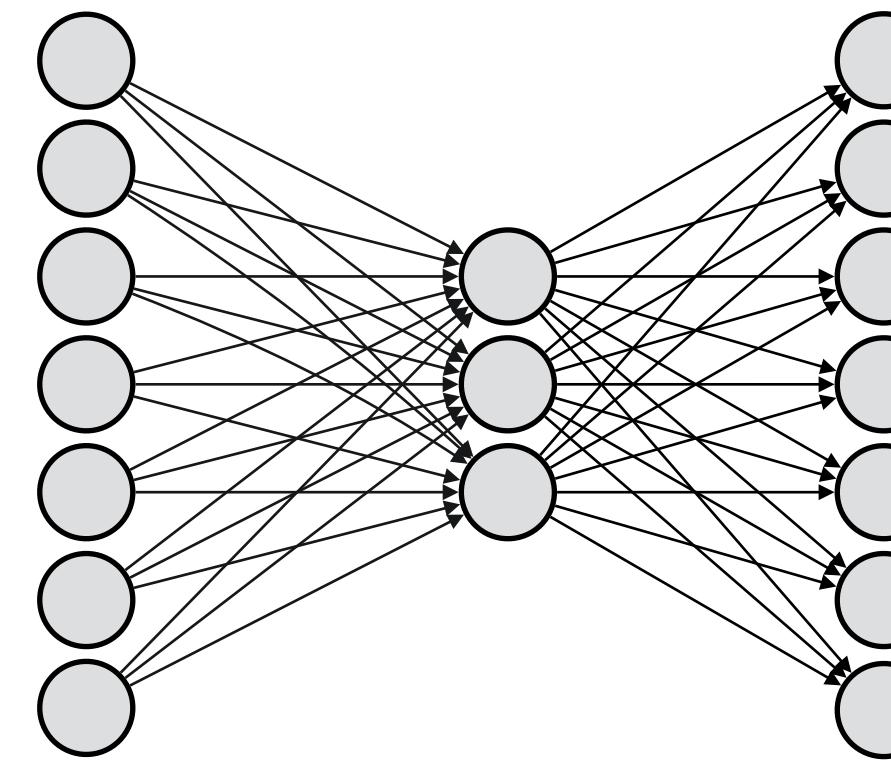
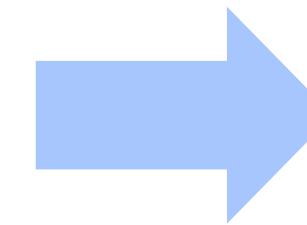
sum over objects (images)

probability of true class

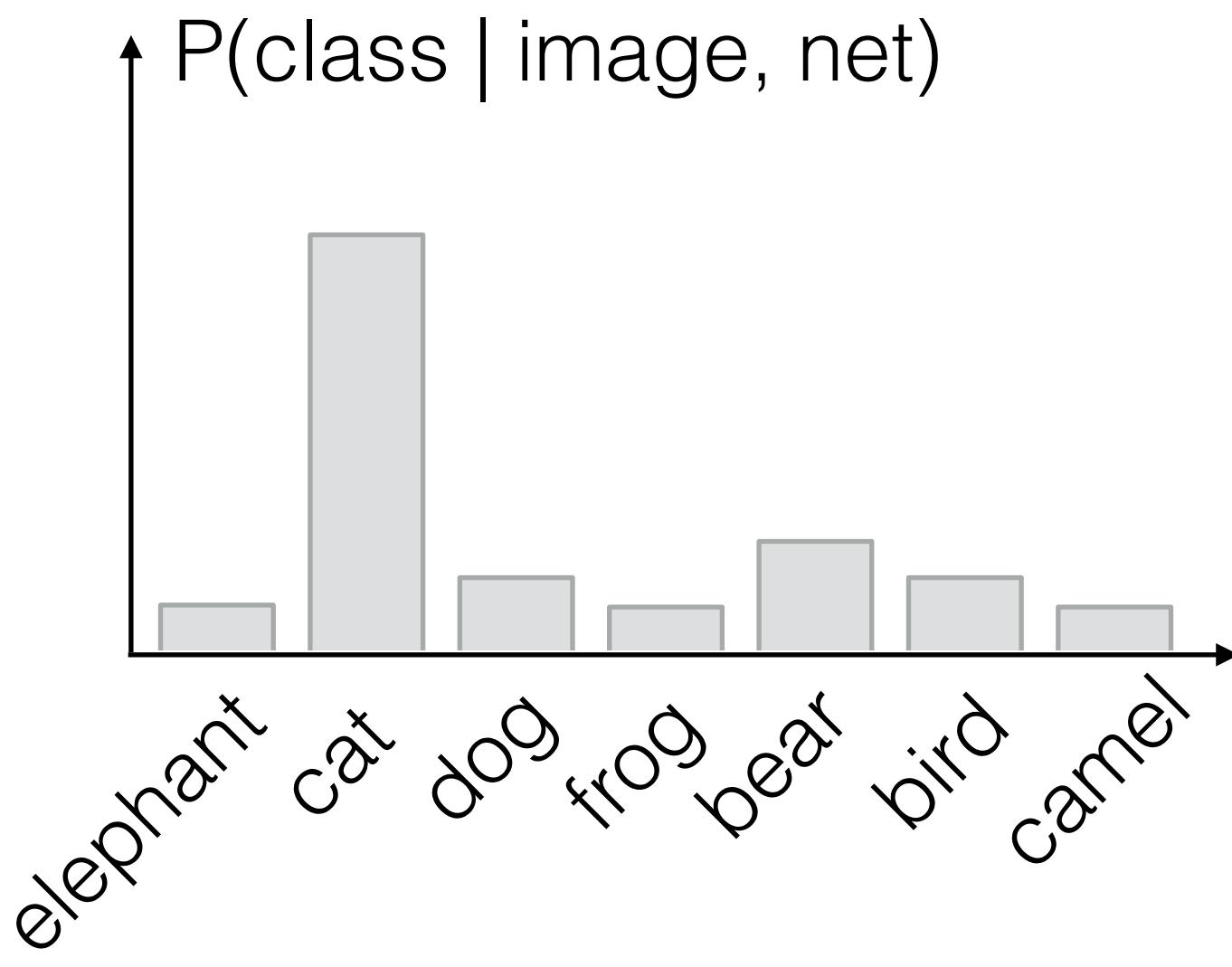
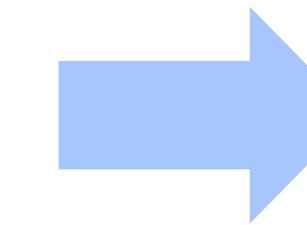
Neural networks



input x



neural network
with weights w



Training — optimization
over weights w
using stochastic
gradient descend:

$$w^{new} = w^{old} + \eta \frac{\partial}{\partial w} \sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, w^{old})$$

$$i_j \sim \text{Unif}(1, \dots, N)$$

m — mini-batch size

η — learning rate

Regularization by noise

- Traditional regularization: add some penalty for model complexity
 - L_2 , L_1 - regularization, max norm constraint

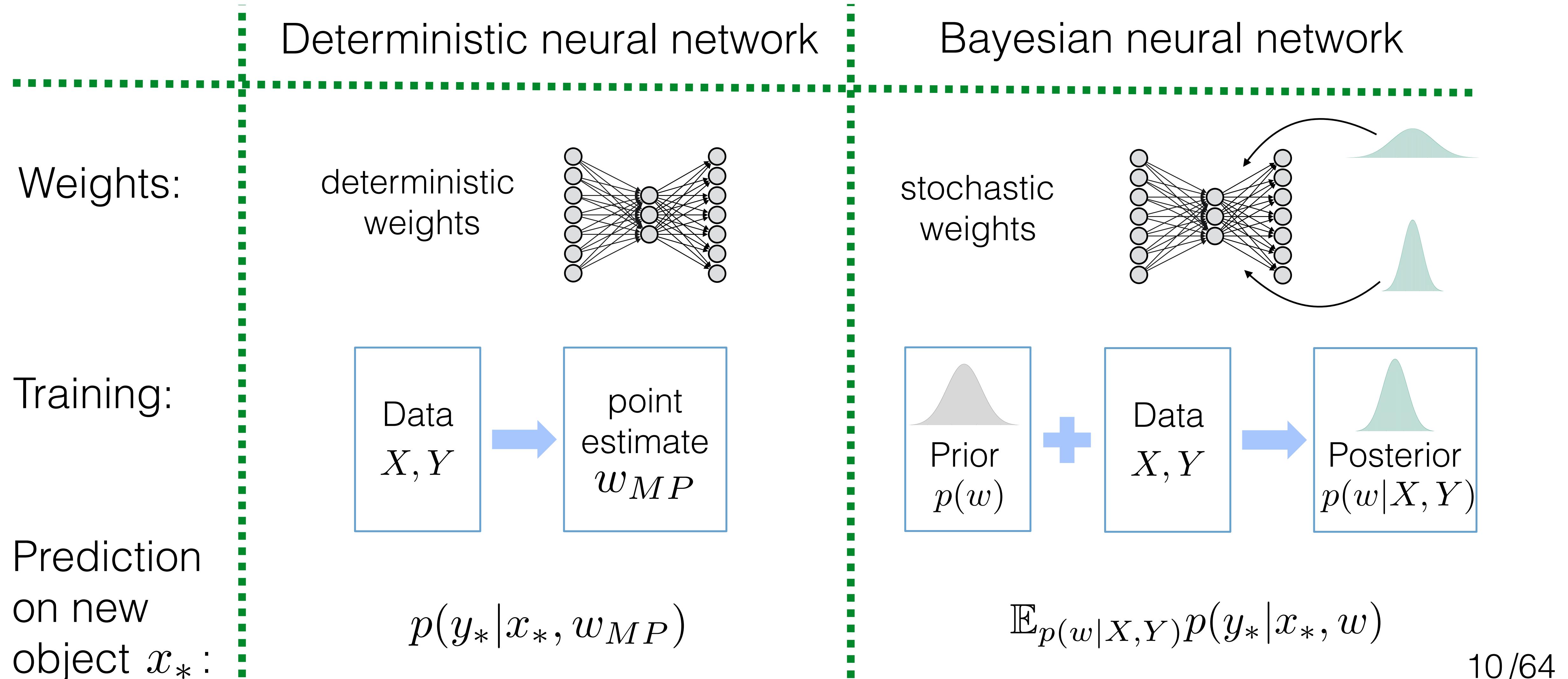
$$\text{Objective} = \text{DataLoss}(X, Y, w) + \text{Regularizer}(w)$$

- More recent approaches: regularization by noise
 - Data augmentation, dropout, gradient noise

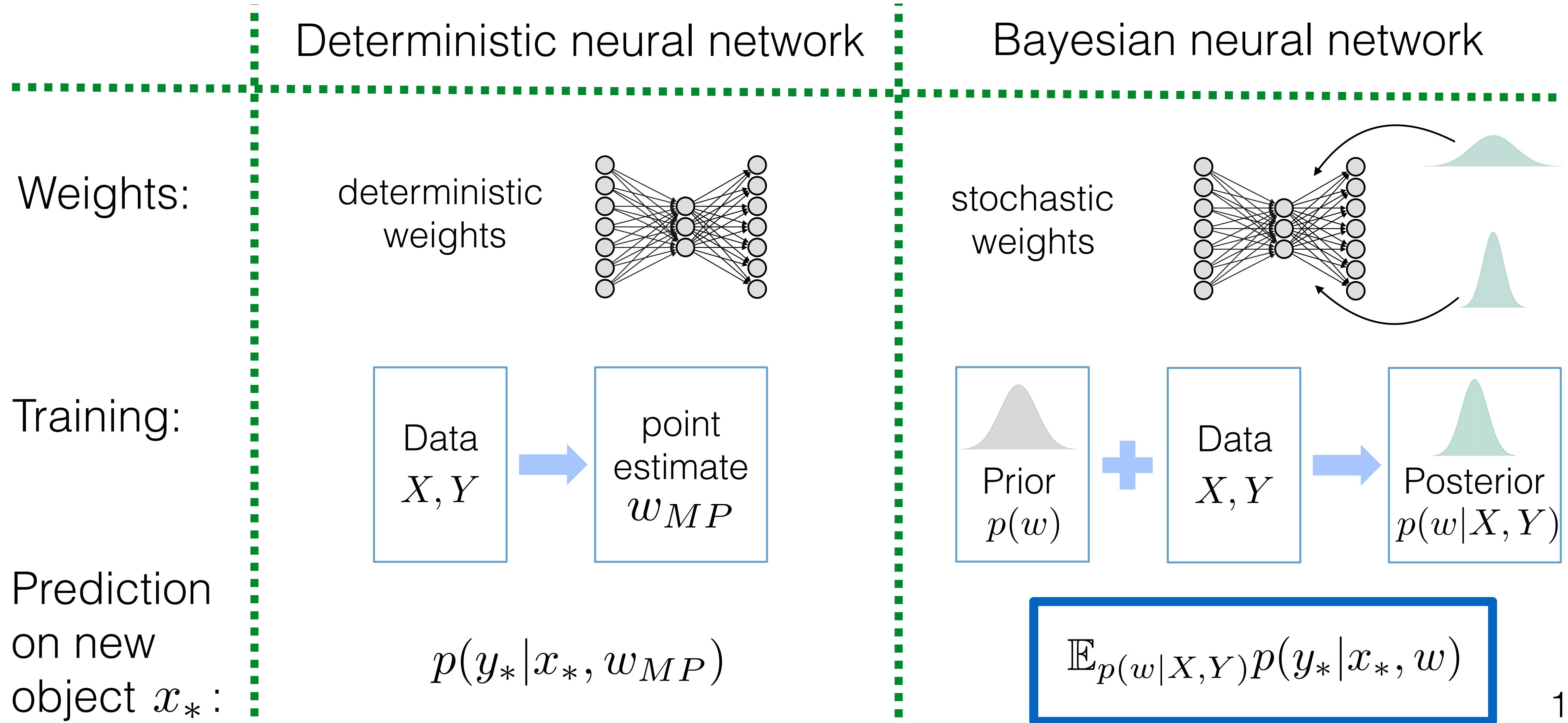
$$\text{Objective} = \mathbb{E}_{p(\Omega)} \text{DataLoss}(X, Y, w, \Omega)$$

Bayesian framework provides a principled approach to training with noise!

Bayesian neural networks



Bayesian neural networks



BNN as an ensemble of neural networks

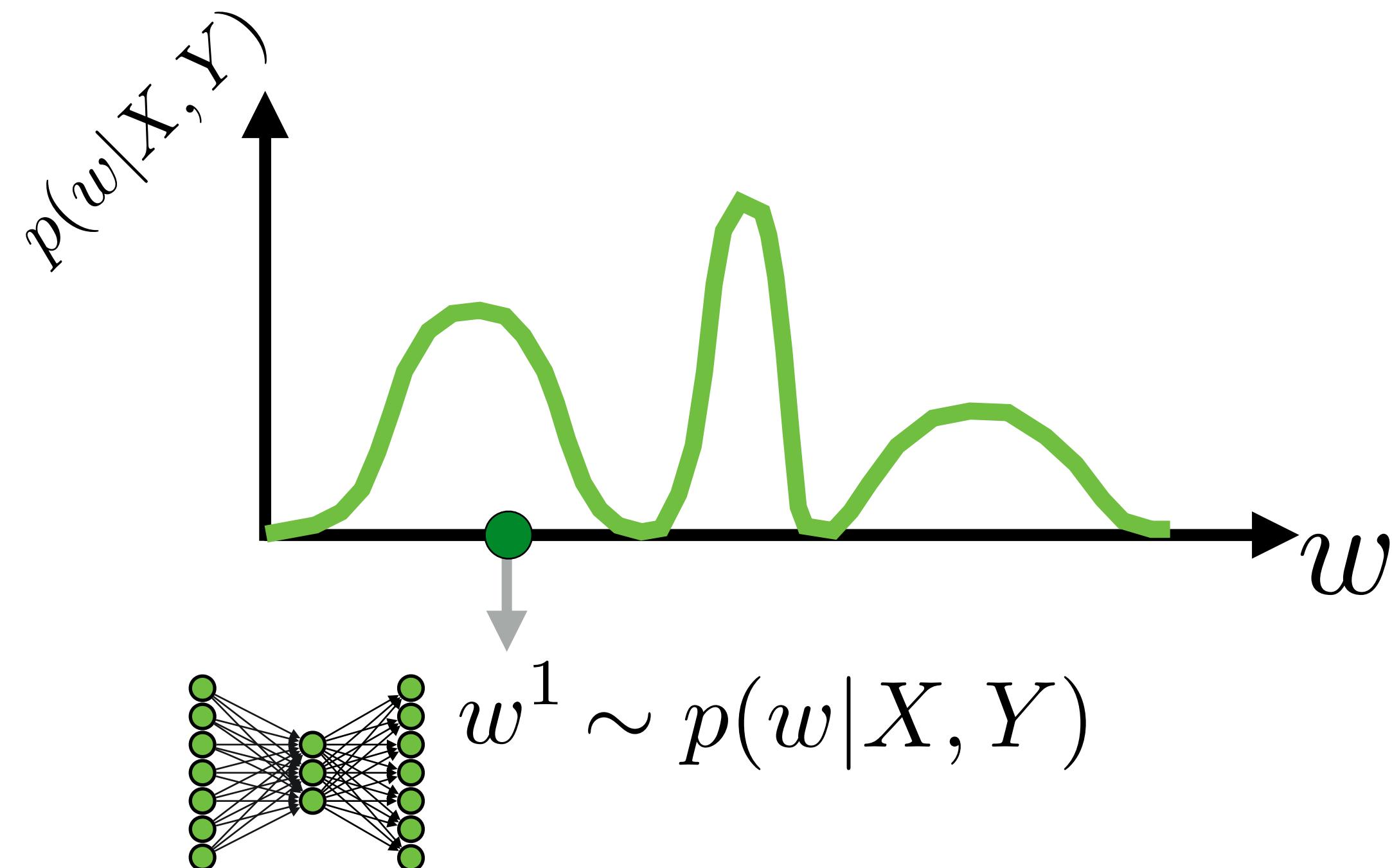
Prediction on a new object x_* :

$$\mathbb{E}_{p(w|X,Y)} p(y_*|x_*, w) \approx \frac{1}{K} \sum_{k=1}^K p(y_*|x_*, w^k), \quad w^k \sim p(w|X, Y)$$

BNN as an ensemble of neural networks

Prediction on a new object x_* :

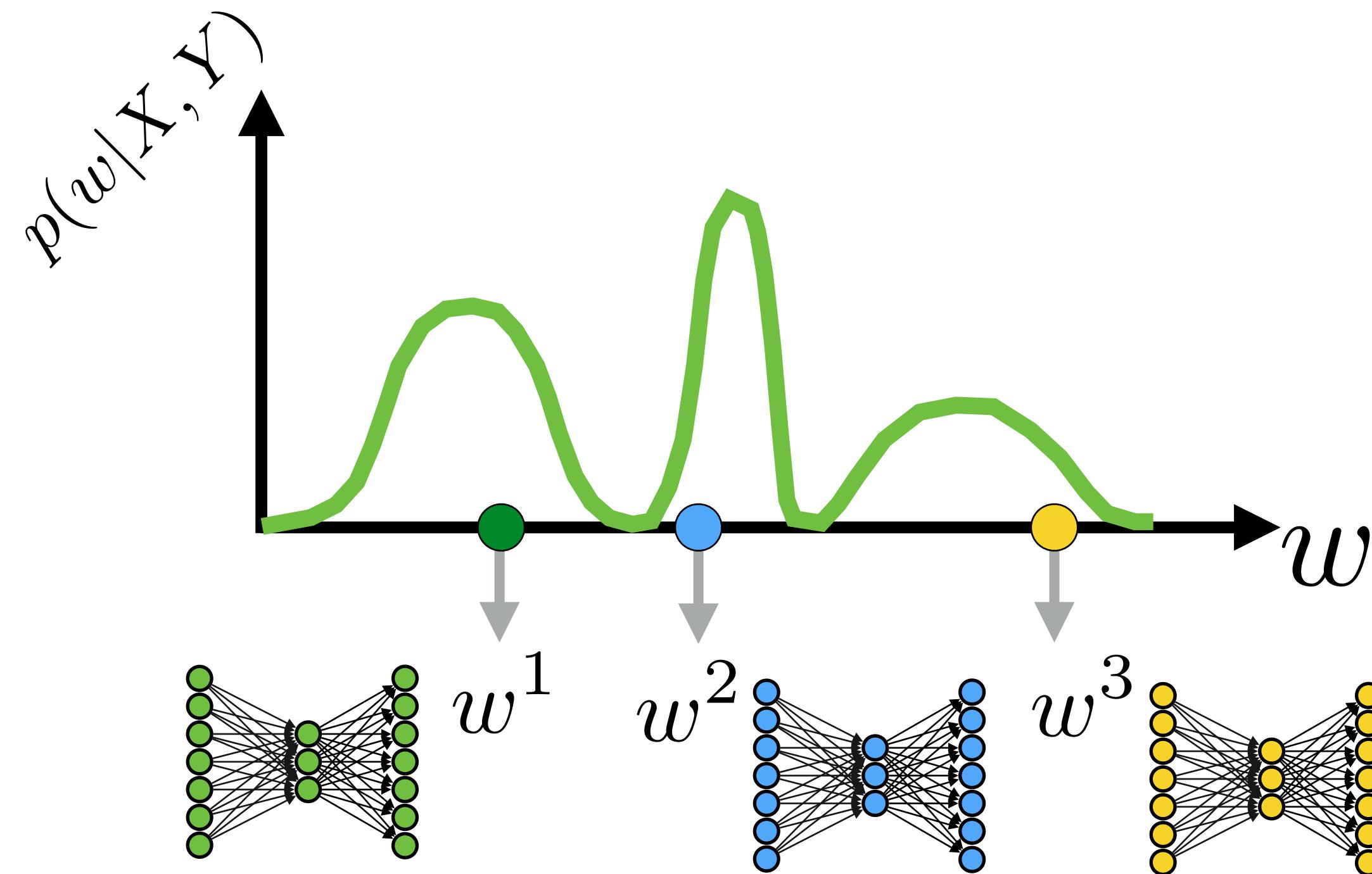
$$\mathbb{E}_{p(w|X,Y)} p(y_*|x_*, w) \approx \frac{1}{K} \sum_{k=1}^K p(y_*|x_*, w^k), \quad w^k \sim p(w|X, Y)$$



BNN as an ensemble of neural networks

Prediction on a new object x_* :

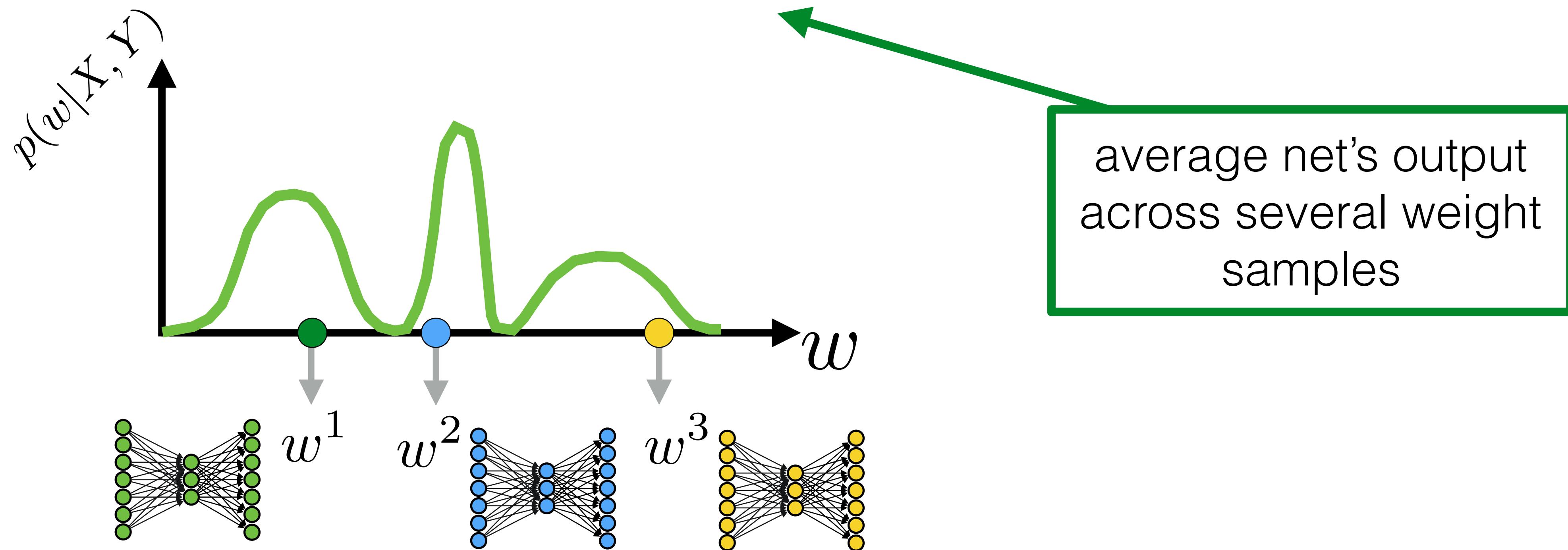
$$\mathbb{E}_{p(w|X,Y)} p(y_*|x_*, w) \approx \frac{1}{K} \sum_{k=1}^K p(y_*|x_*, w^k), \quad w^k \sim p(w|X, Y)$$



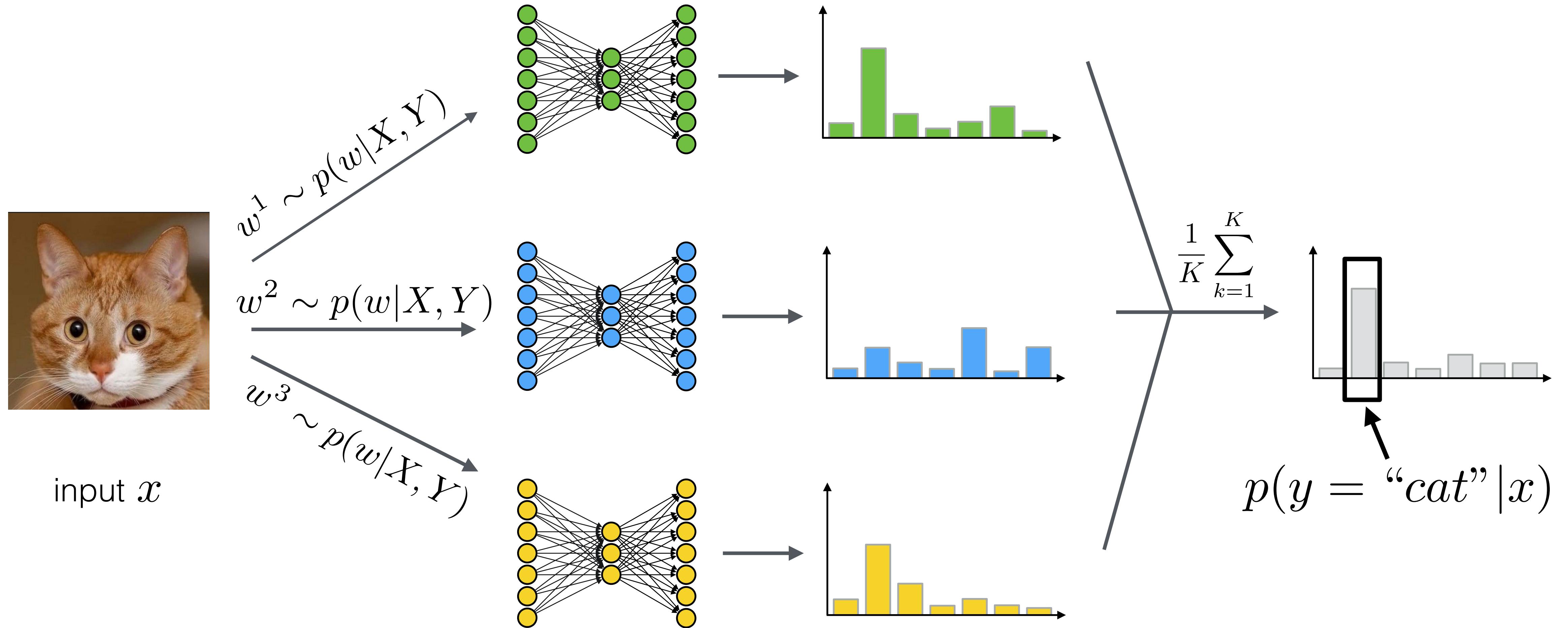
BNN as an ensemble of neural networks

Prediction on a new object x_* :

$$\mathbb{E}_{p(w|X,Y)} p(y_*|x_*, w) \approx \frac{1}{K} \sum_{k=1}^K p(y_*|x_*, w^k), \quad w^k \sim p(w|X, Y)$$



BNN as an ensemble of neural networks



BNN as an ensemble of neural networks

Prediction on a new object x_* :

$$\mathbb{E}_{p(w|X,Y)} p(y_*|x_*, w) \approx \frac{1}{K} \sum_{k=1}^K p(y_*|x_*, w^k), \quad w^k \sim p(w|X, Y)$$

- Higher quality (models compensate each other's errors)
- Better uncertainty estimation

Questions

Questions

- How many forward passes do we perform to make the prediction in the BNN?

Questions

- How many forward passes do we perform to make the prediction in the BNN?
- If the posterior is a delta function: $p(w|X,Y) = \delta(w_0)$,
how will the predictive ensemble look like?

Questions

- How many forward passes do we perform to make the prediction in the BNN?
- If the posterior is a delta function: $p(w|X,Y) = \delta(w_0)$,
how will the predictive ensemble look like?

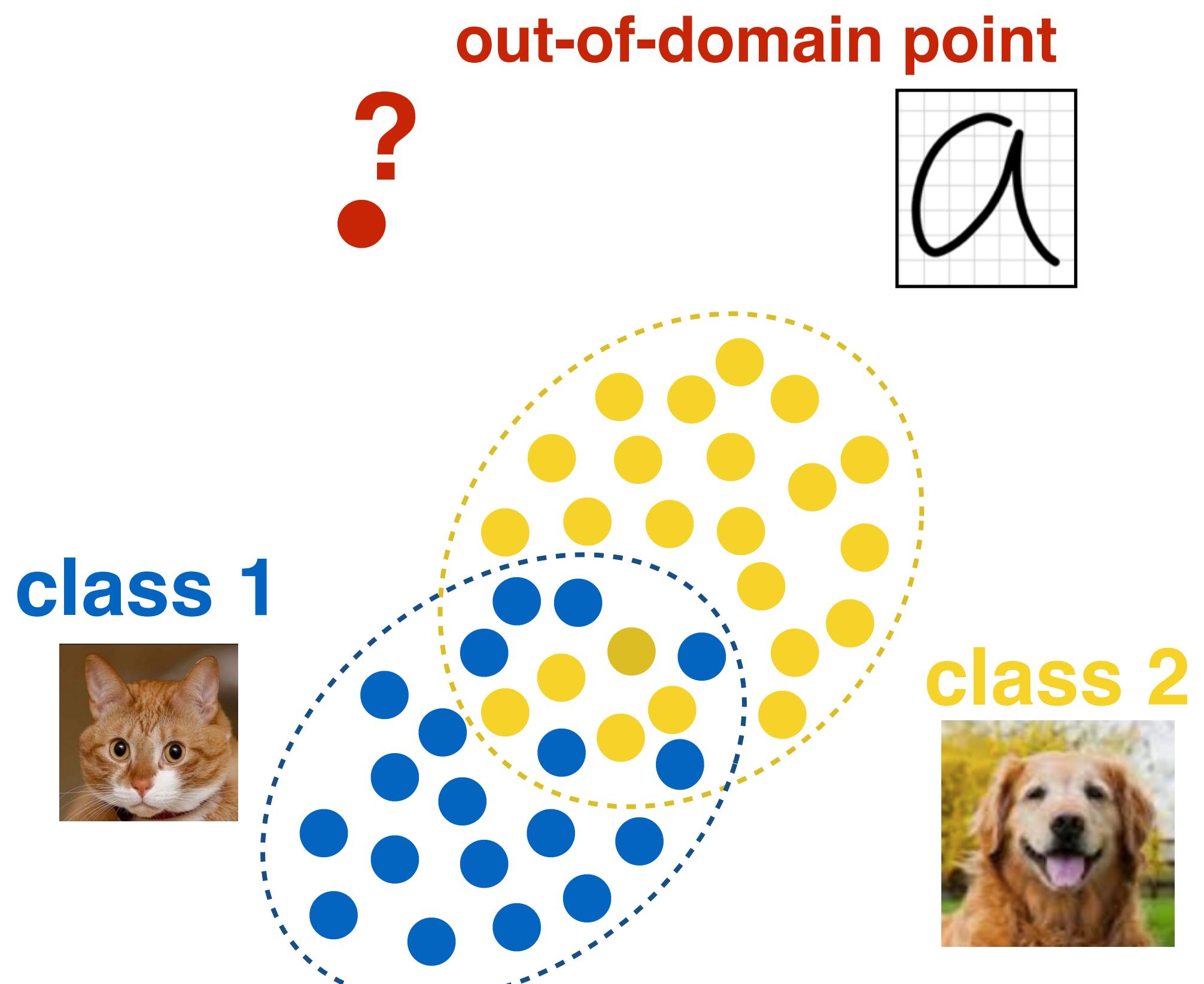
$w \sim p(w|X,Y) = \delta(w_0) \Leftrightarrow w \equiv w_0 \longrightarrow$ all the samples
(networks) are
the same

Why go Bayesian?

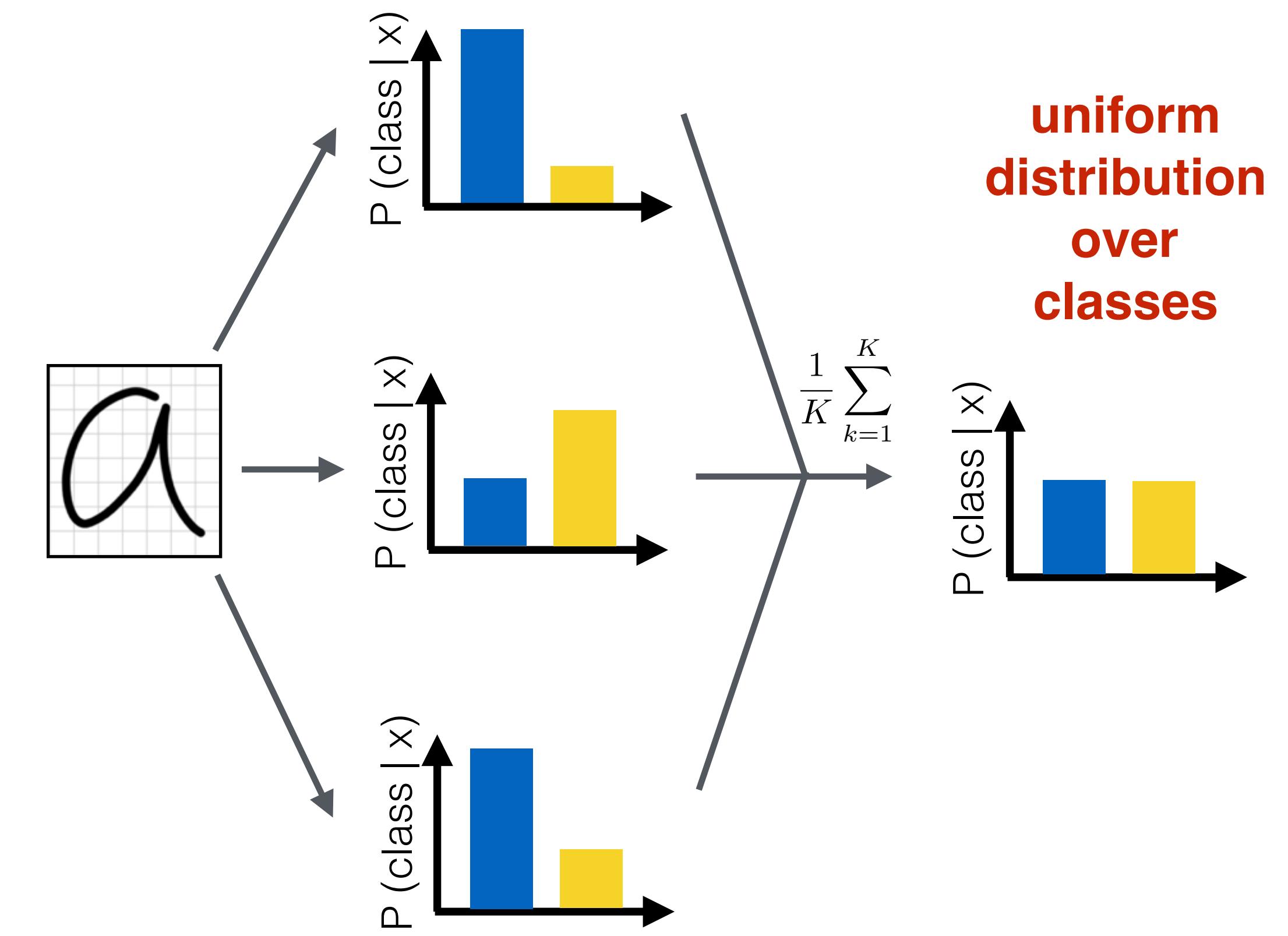
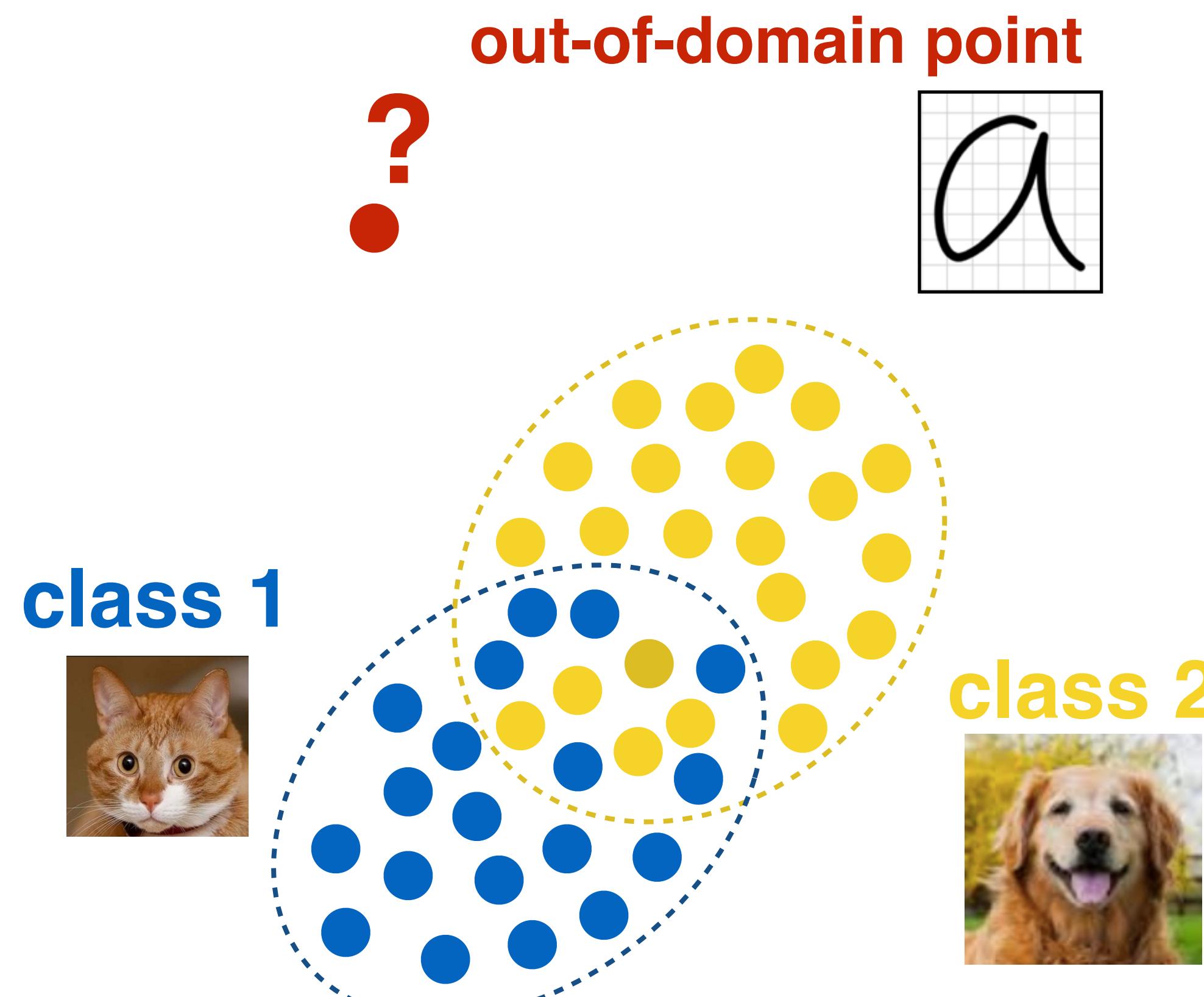
A principled framework with many useful applications

- Regularization
- Ensembling
- Uncertainty estimation
- On-line / continual learning
- Different priors lead to different properties of the network
- Automatic hyperparameter choice

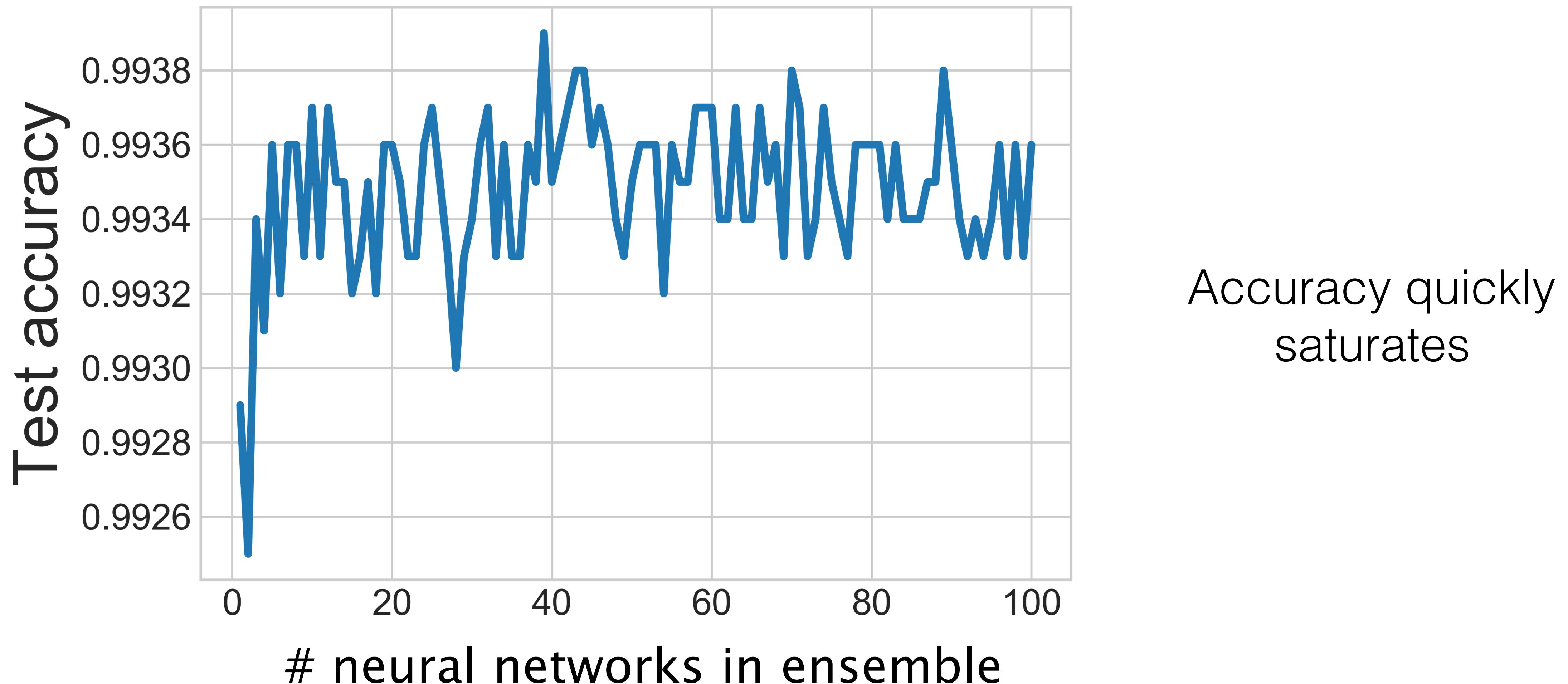
Uncertainty estimation



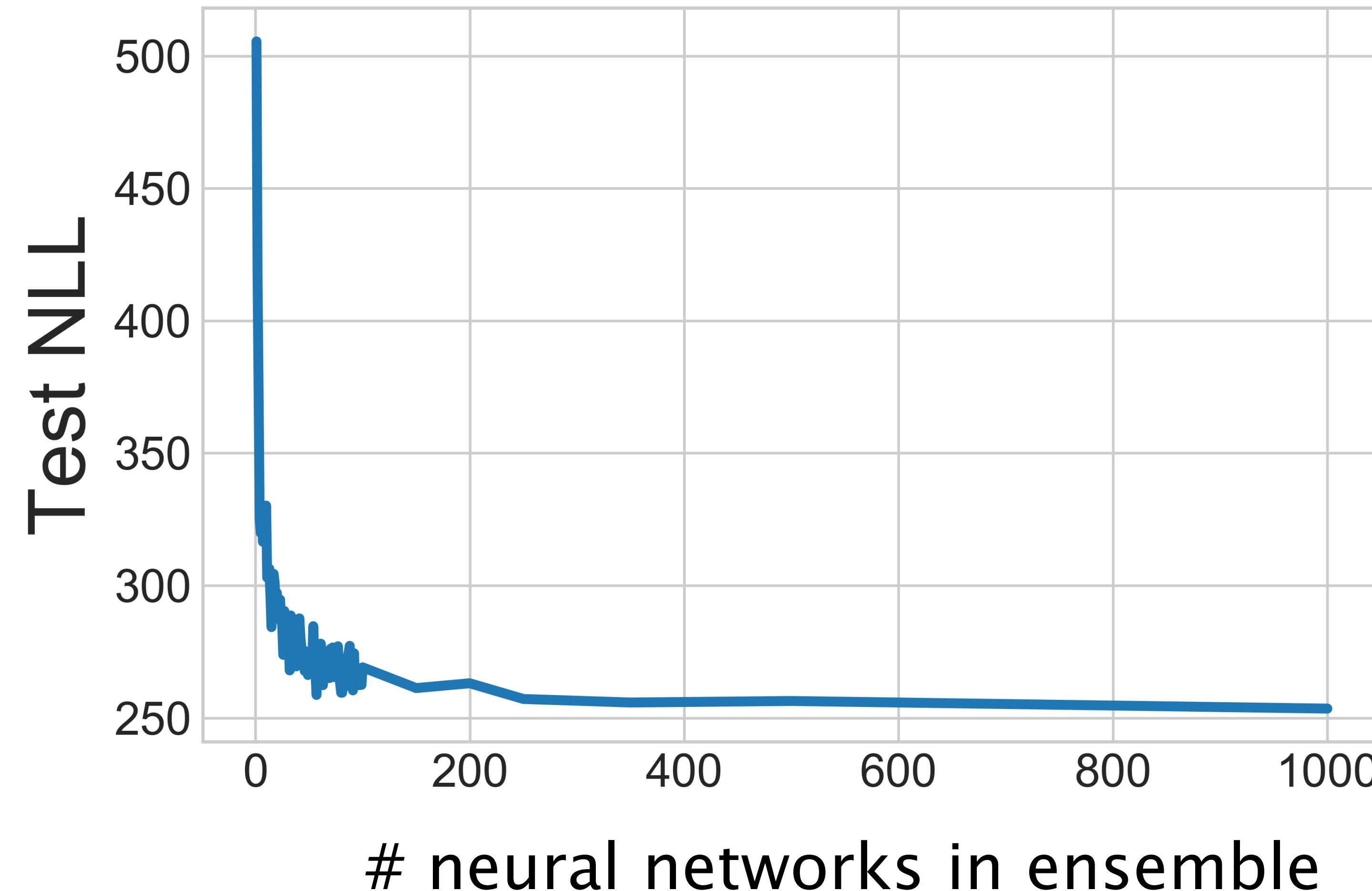
Uncertainty estimation



Ensembling: accuracy



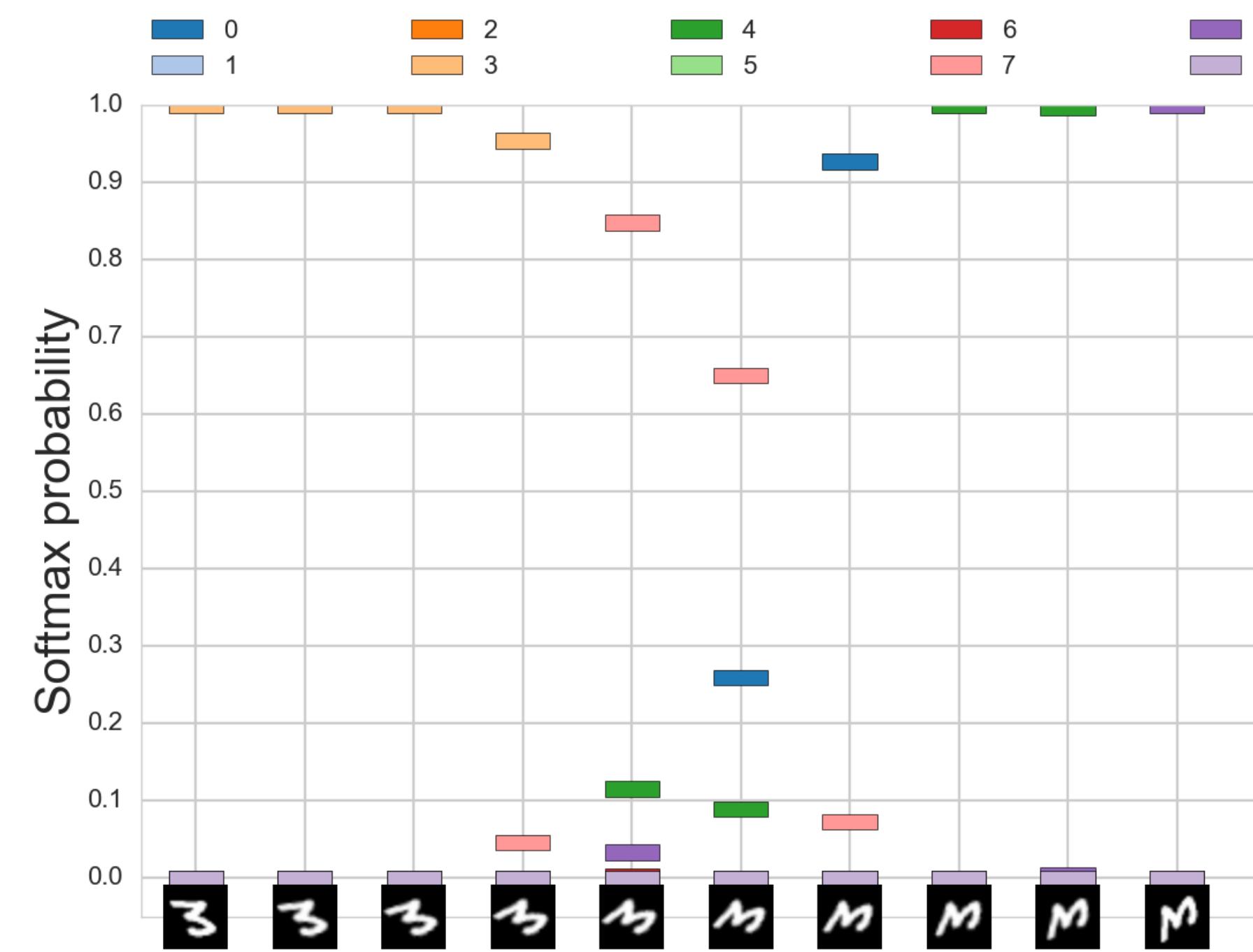
Ensembling: uncertainty estimation



But the negative
log—likelihood
keeps improving!
This is a measure
of “**uncertainty**”

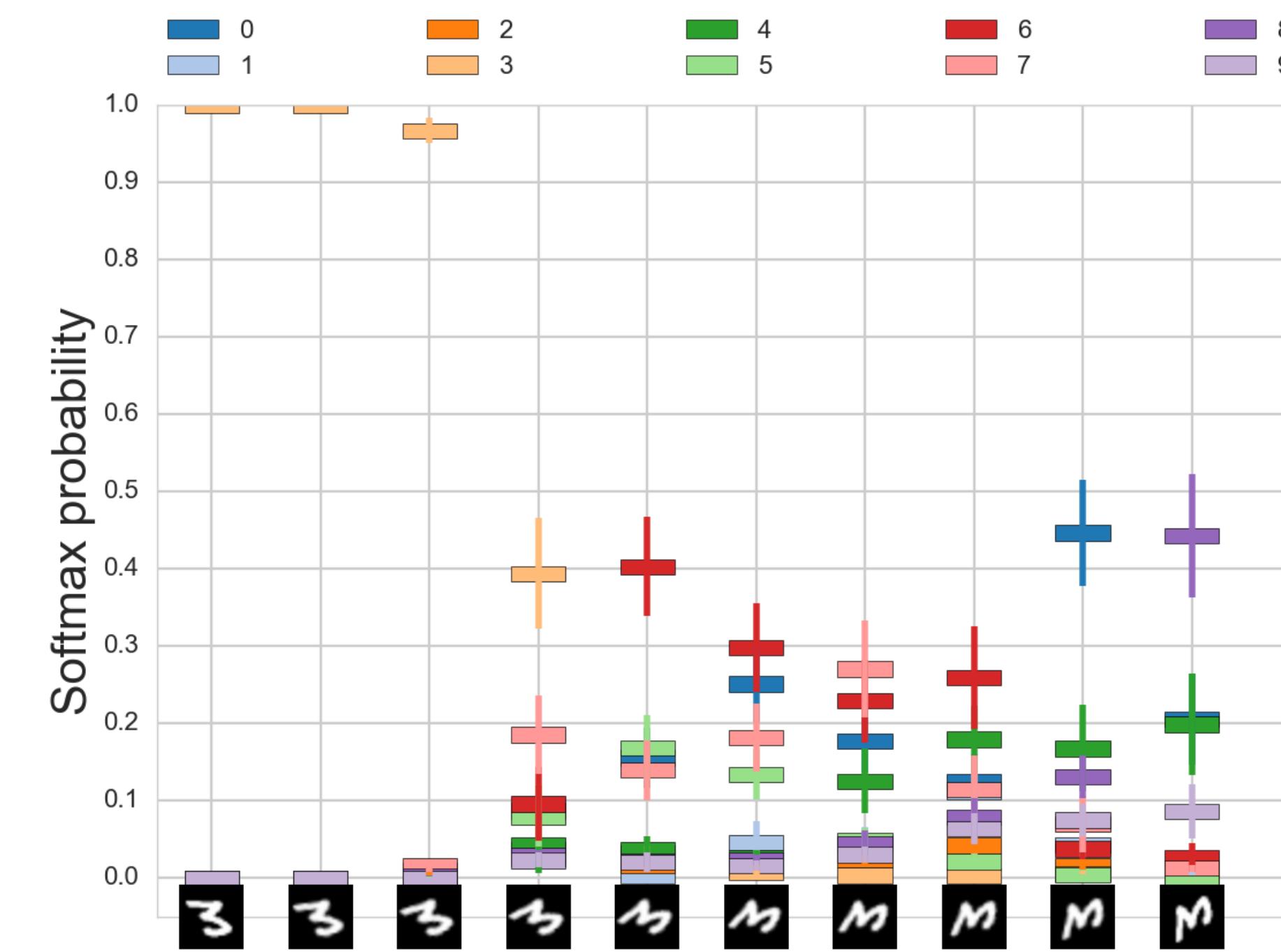
Uncertainty in classification: experiment

Deterministic NN



(a) LeNet with weight decay

Bayesian NN



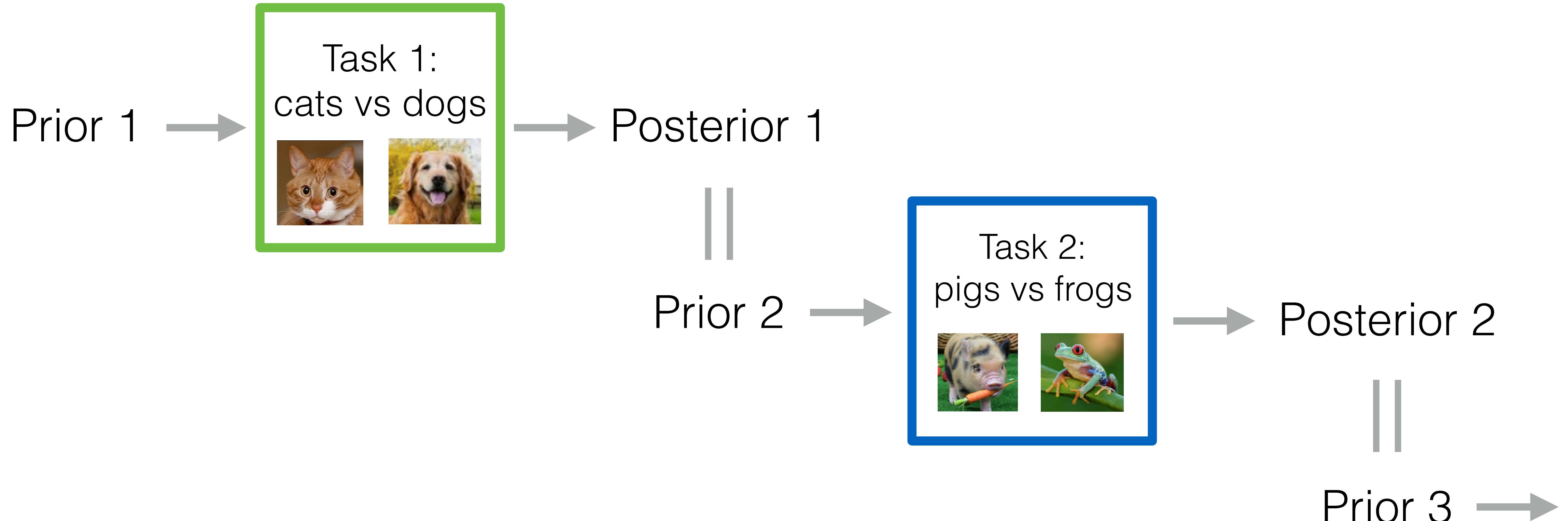
(b) LeNet with multiplicative formalizing flows

Why go Bayesian?

A principled framework with many useful applications

- Regularization ✓
- Ensembling ✓
- Uncertainty estimation ✓
- On-line / continual learning
- Different priors lead to different properties of the network
- Automatic hyperparameter choice

On-line / continual learning

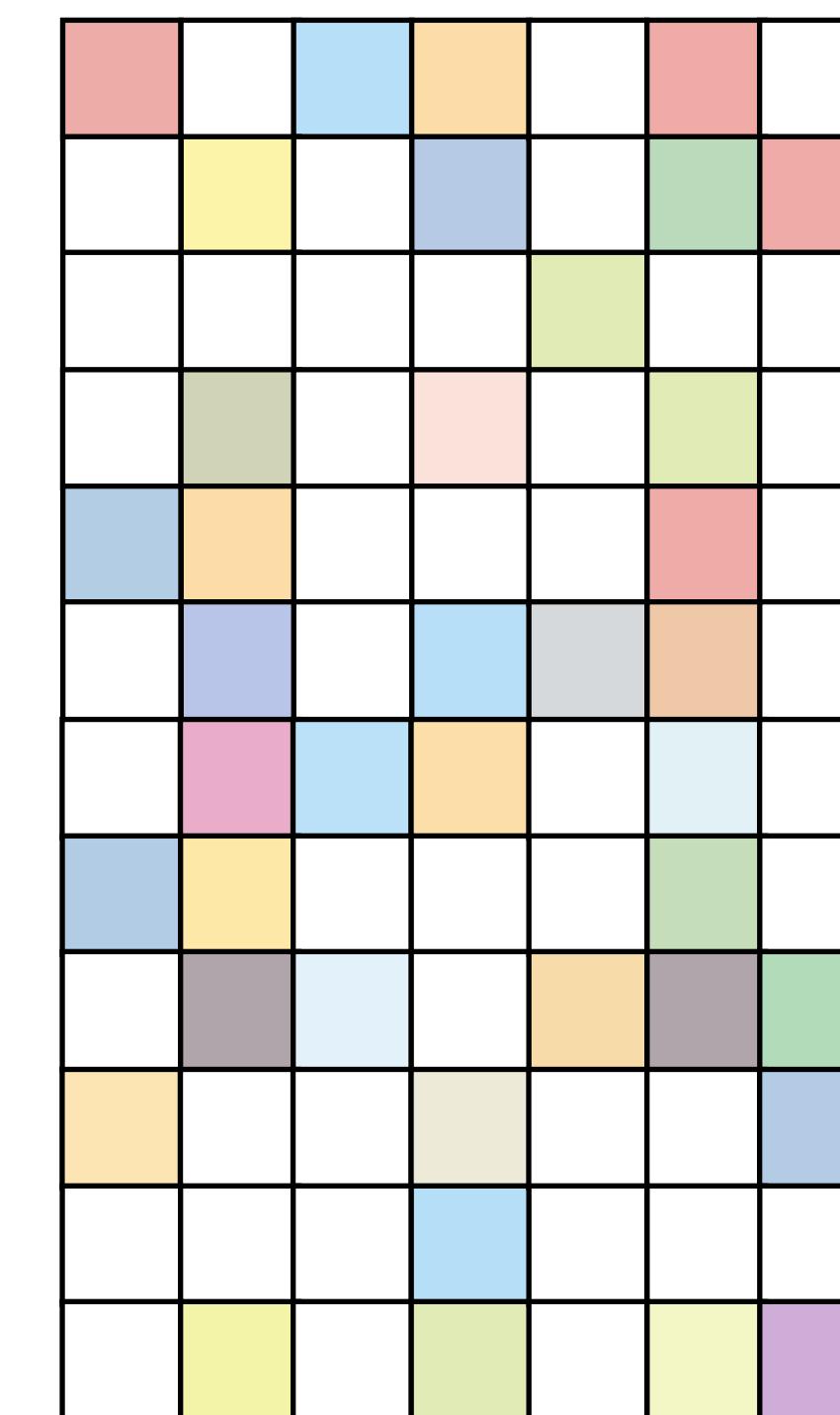
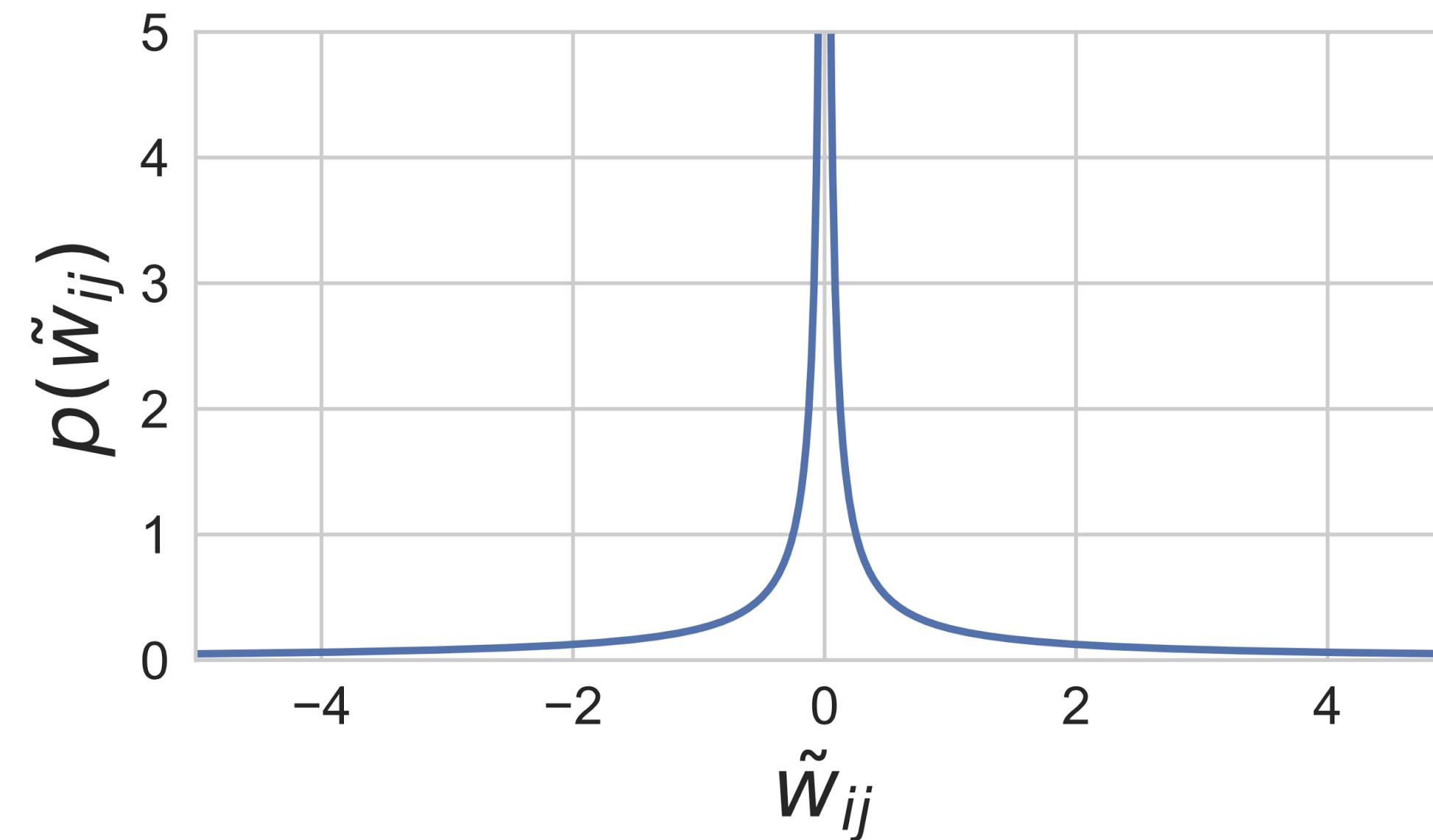


Prior can encode our desirable model properties

Prior concentrated at zero



A lot of zero weights



Weight matrix W

Why go Bayesian?

A principled framework with many useful applications

- Regularization ✓
- Ensembling ✓
- Uncertainty estimation ✓
- On-line / continual learning ✓
- Different priors lead to different properties of the network ✓
- Automatic hyperparameter choice

Plan

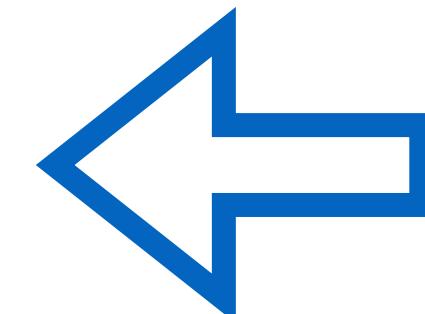
- Advantages of using Bayesian neural networks
- Training Bayesian neural networks
- Q&A + exercises

Training methods: summary

Probabilistic model: $p(Y|X, w)p(w)$

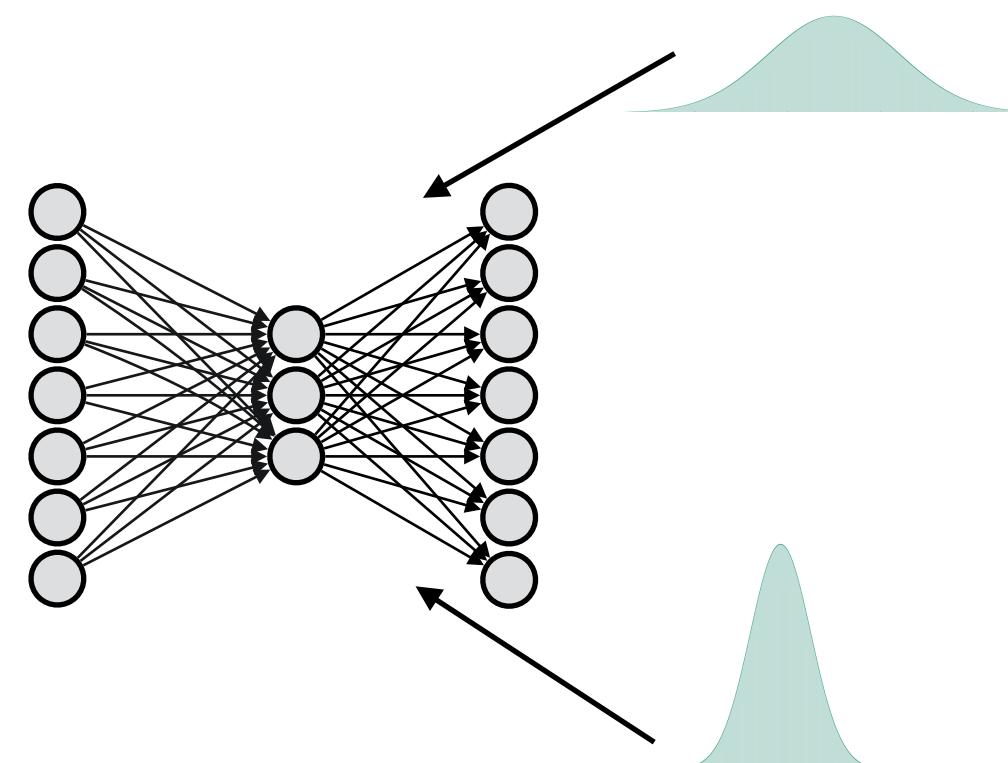
We want to compute: $p(w|X, Y)$

Approximation		Inference
Exact	$p(w X, Y)$	Full Bayesian inference
Parametric	$p(w X, Y) \approx q(w \lambda)$	Parametric Var. Inference
Delta function	$p(w X, Y) \approx \delta(w_{MP})$	Max. posterior inference
No prior	w_{ML}	Max. likelihood inference

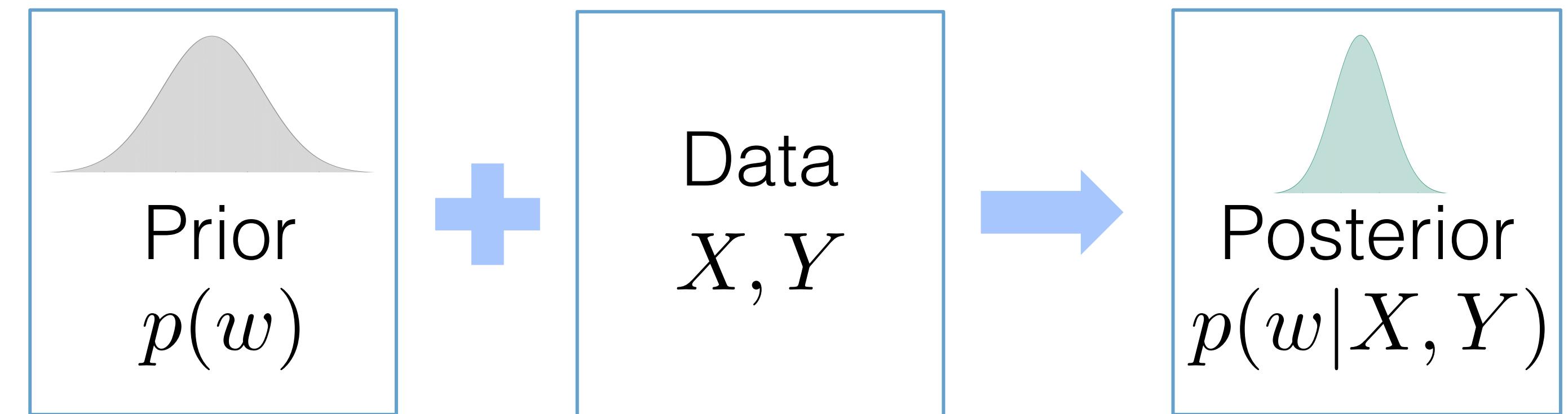


Training Bayesian neural networks

Stochastic weights:



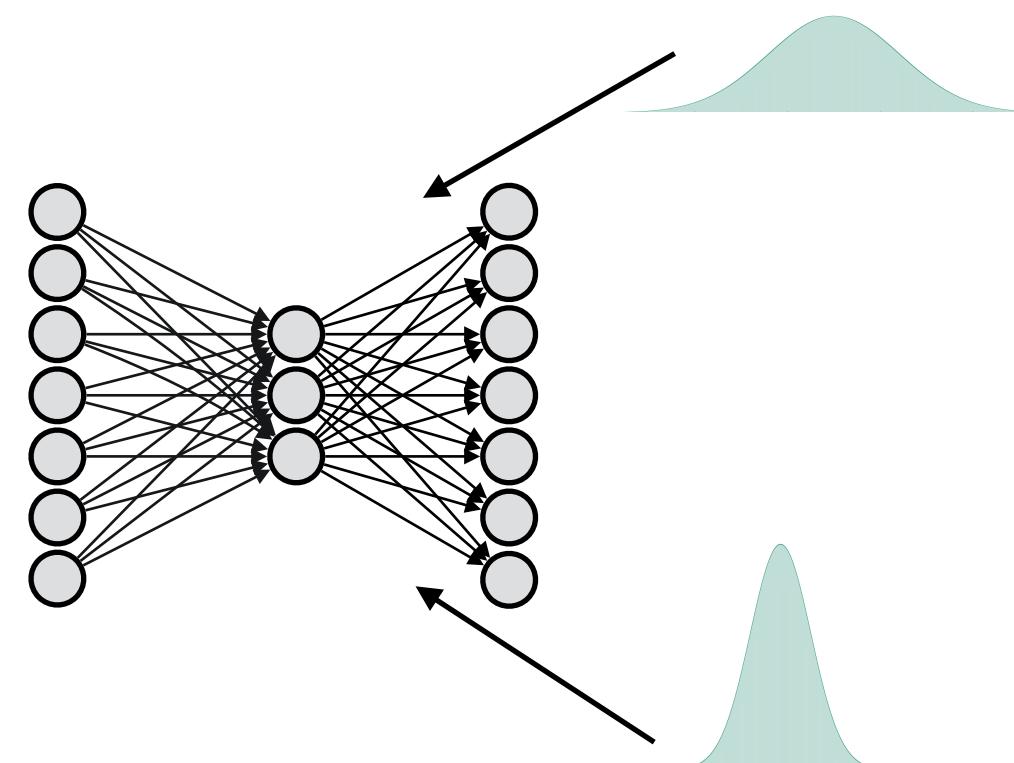
Bayesian Inference:



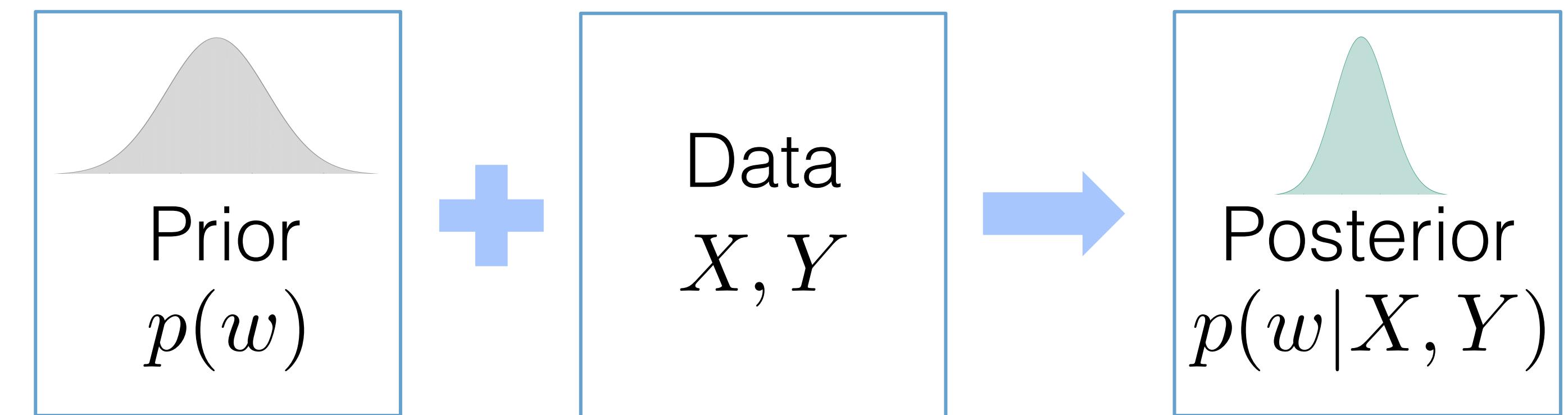
$$p(w|X, Y) = \frac{p(Y|X, w)p(w)}{\int p(Y|X, \tilde{w})p(\tilde{w})d\tilde{w}}$$

Training Bayesian neural networks

Stochastic weights:



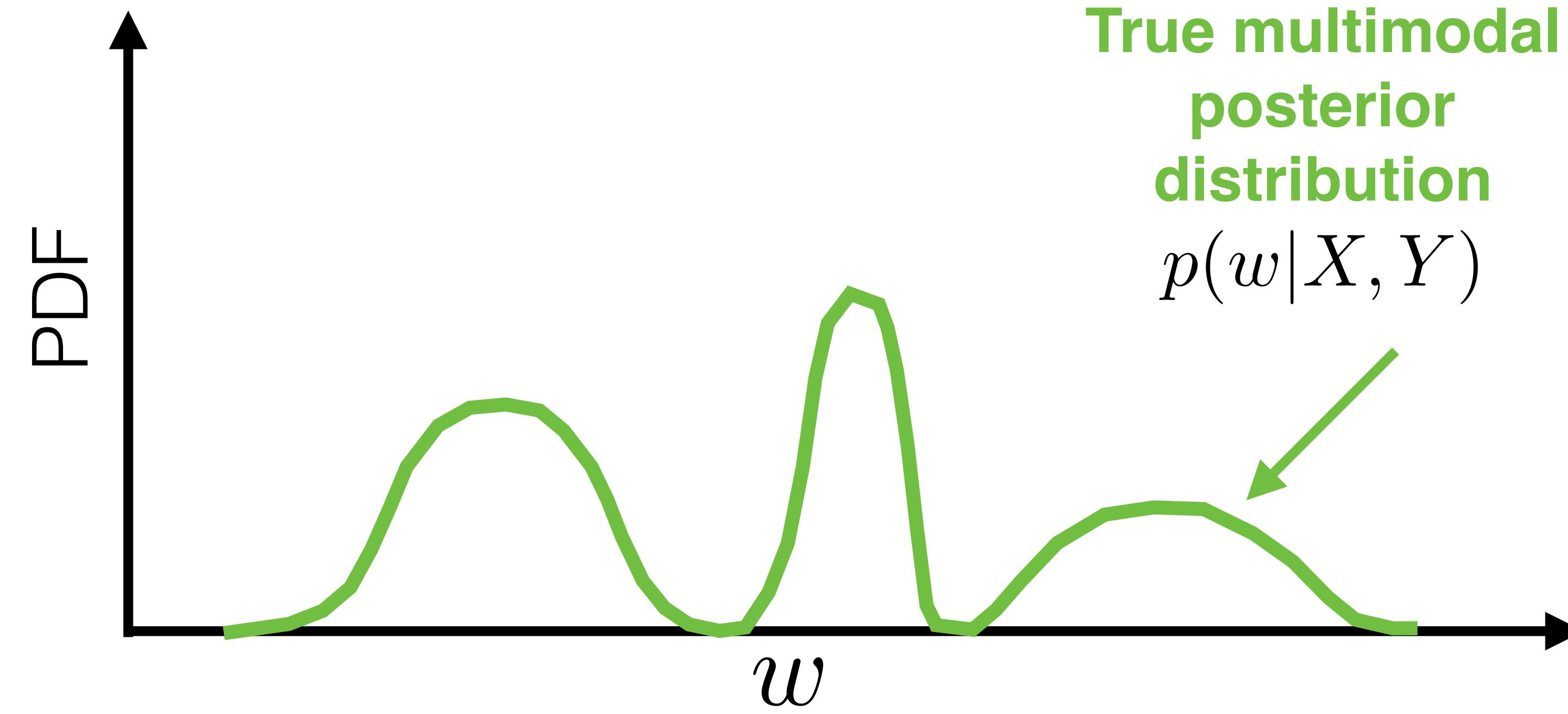
Bayesian Inference:



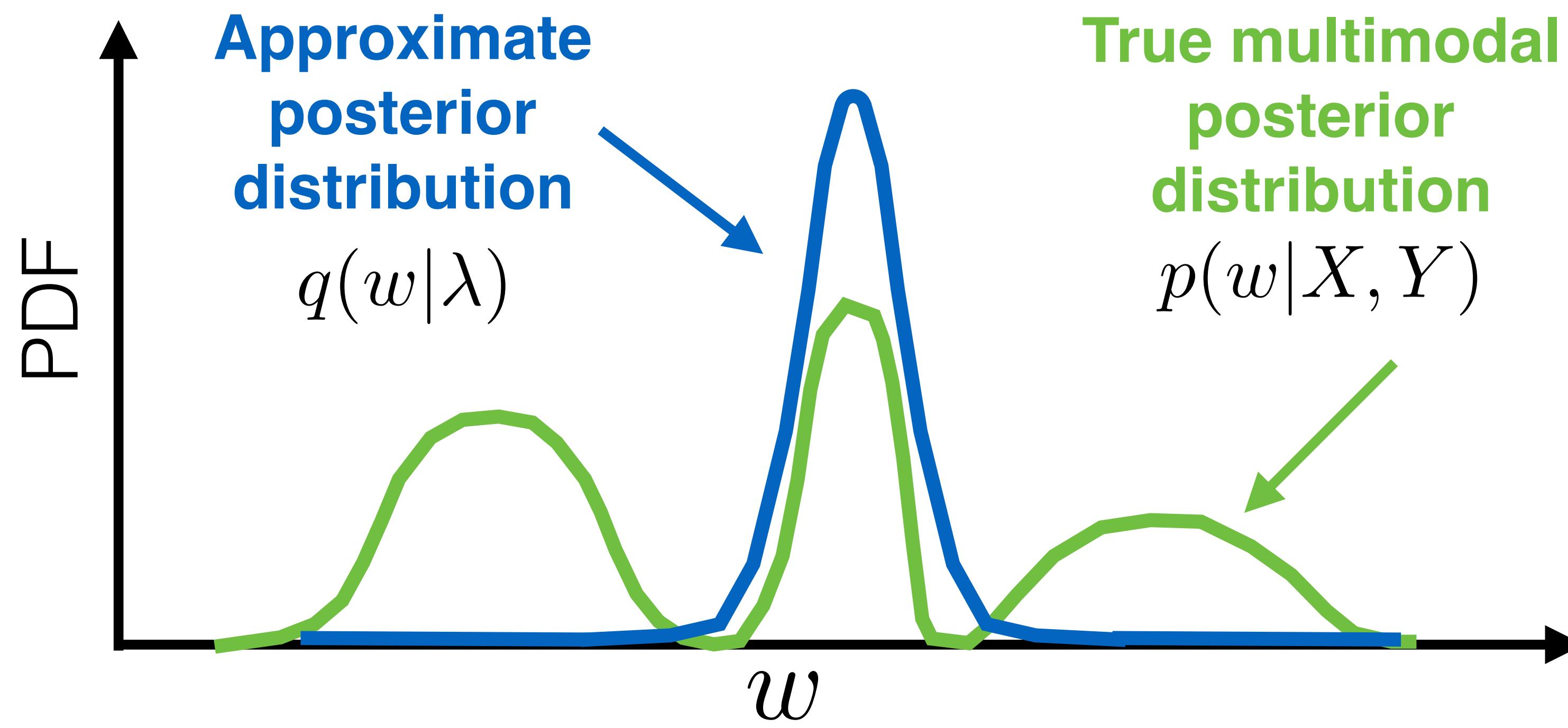
$$p(w|X, Y) = \frac{p(Y|X, w)p(w)}{\int p(Y|X, \tilde{w})p(\tilde{w})d\tilde{w}}$$

Intractable!

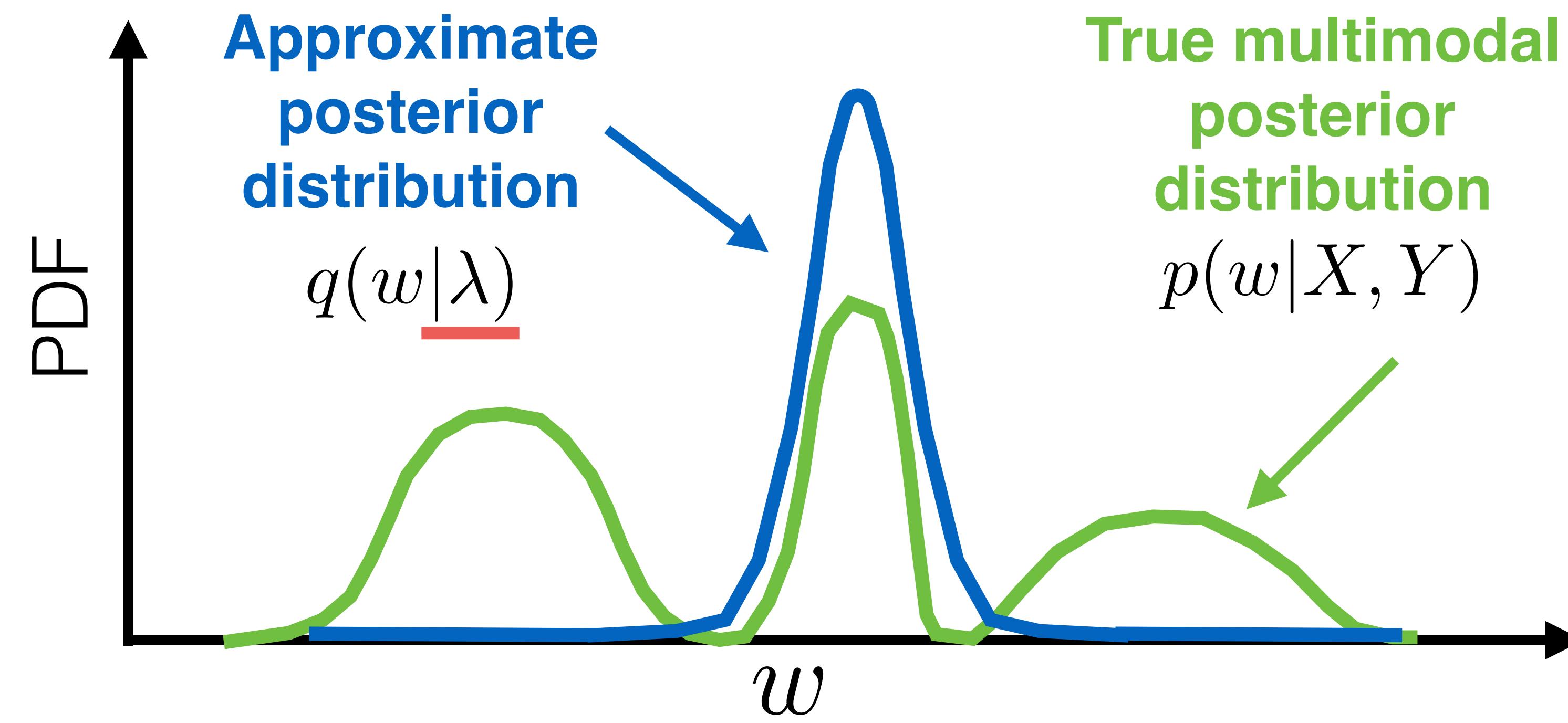
Training Bayesian neural networks



Training Bayesian neural networks



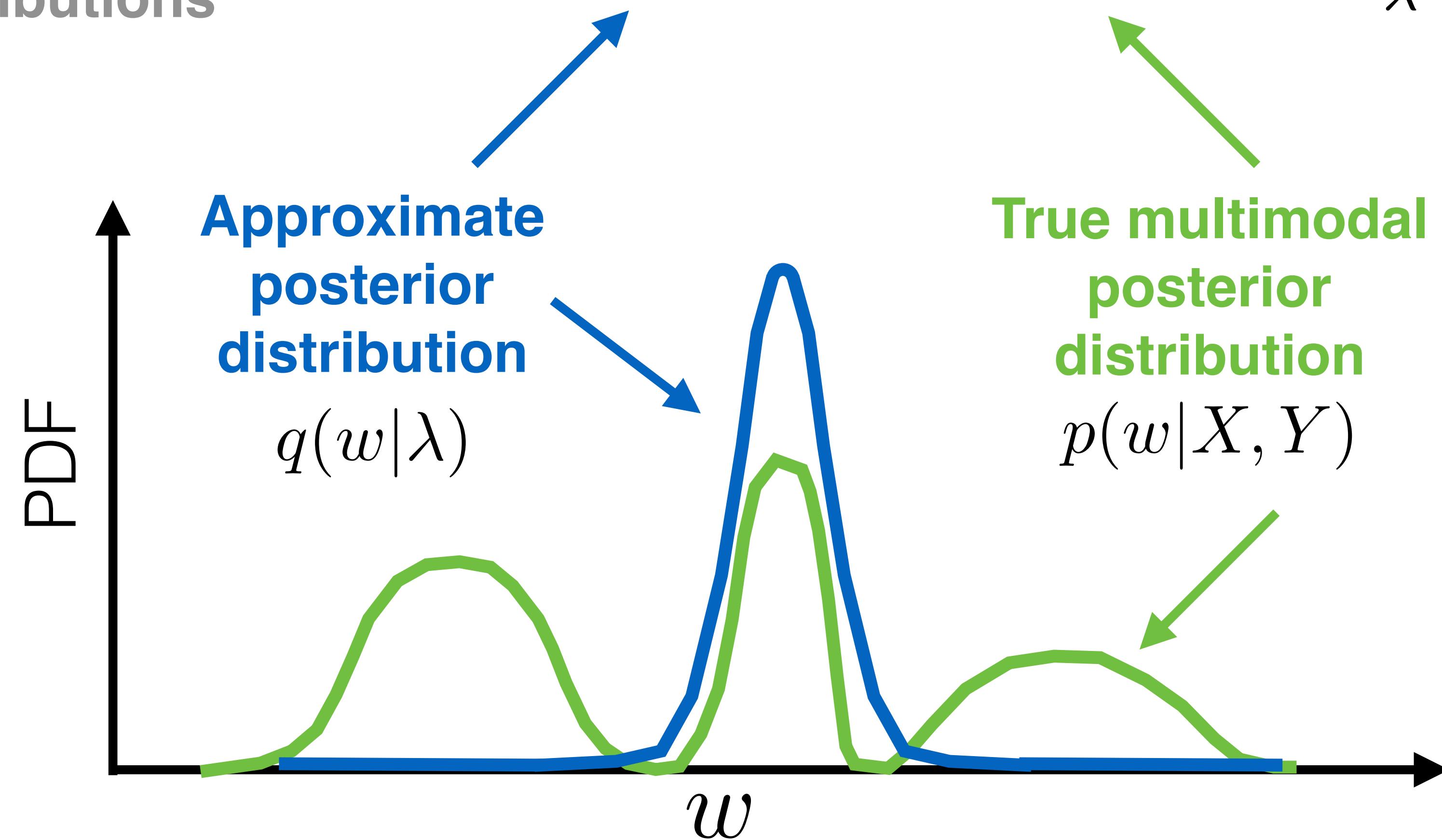
Training Bayesian neural networks



Training Bayesian neural networks

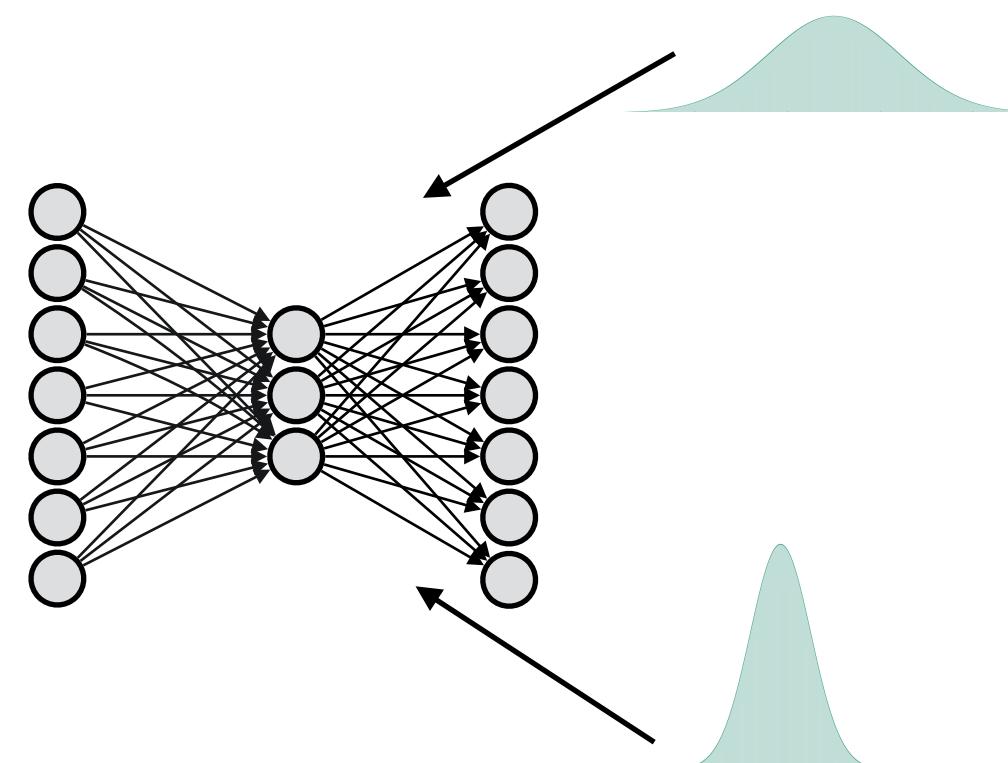
“distance”
between two
distributions

$$\rightarrow KL(q(w|\lambda) || p(w|X, Y)) \rightarrow \min_{\lambda}$$

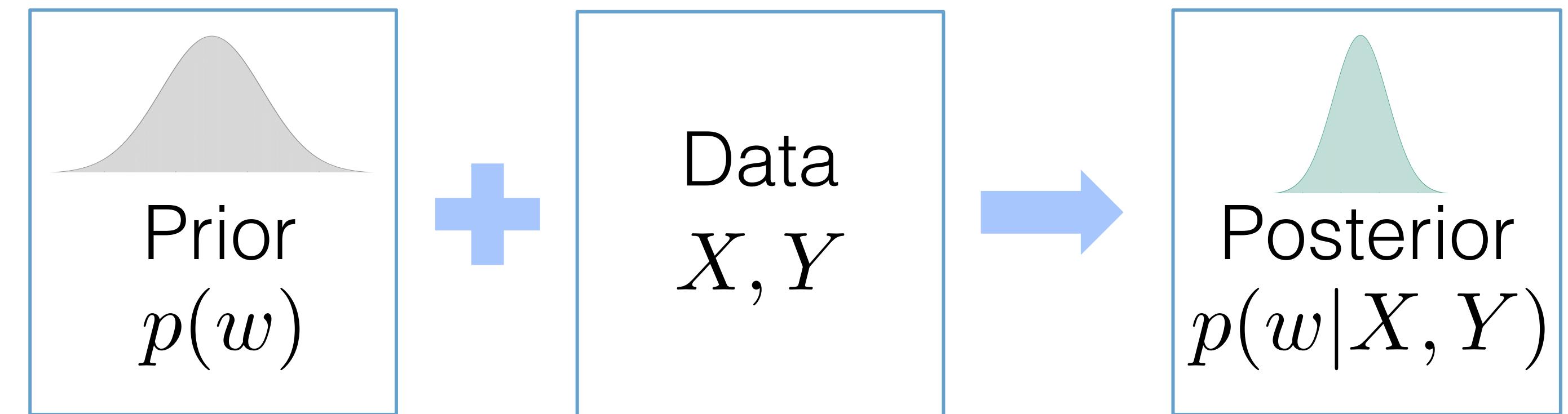


Training Bayesian neural networks

Stochastic weights:



Bayesian Inference:

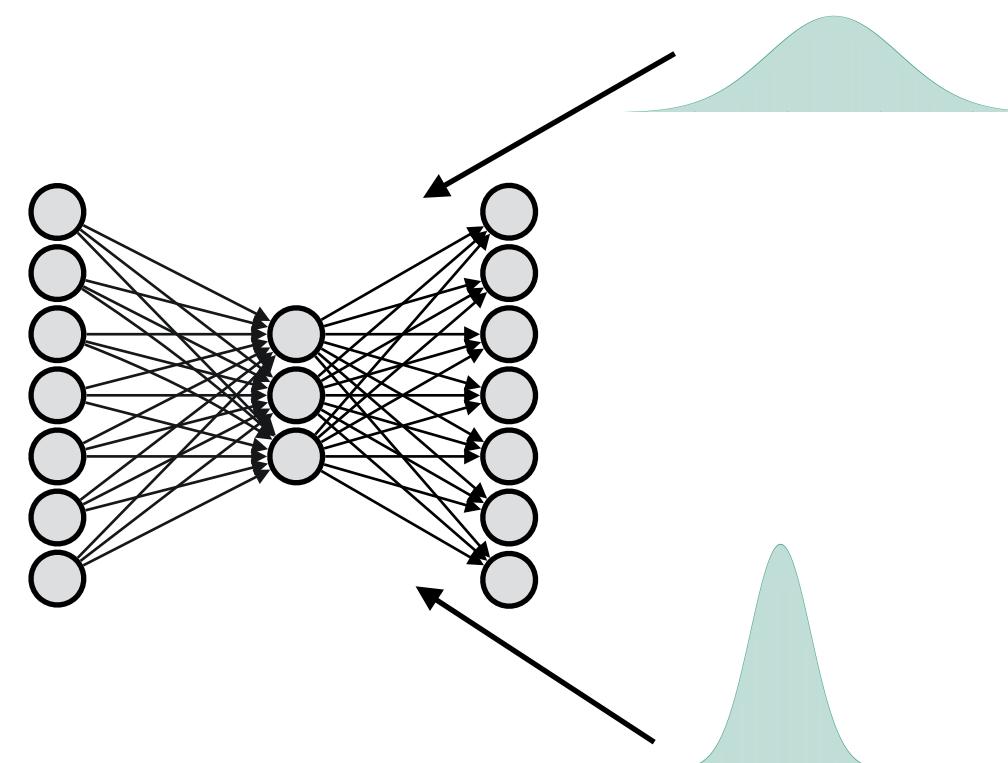


Posterior is intractable in neural networks → approximate it with $q(w|\lambda)$:

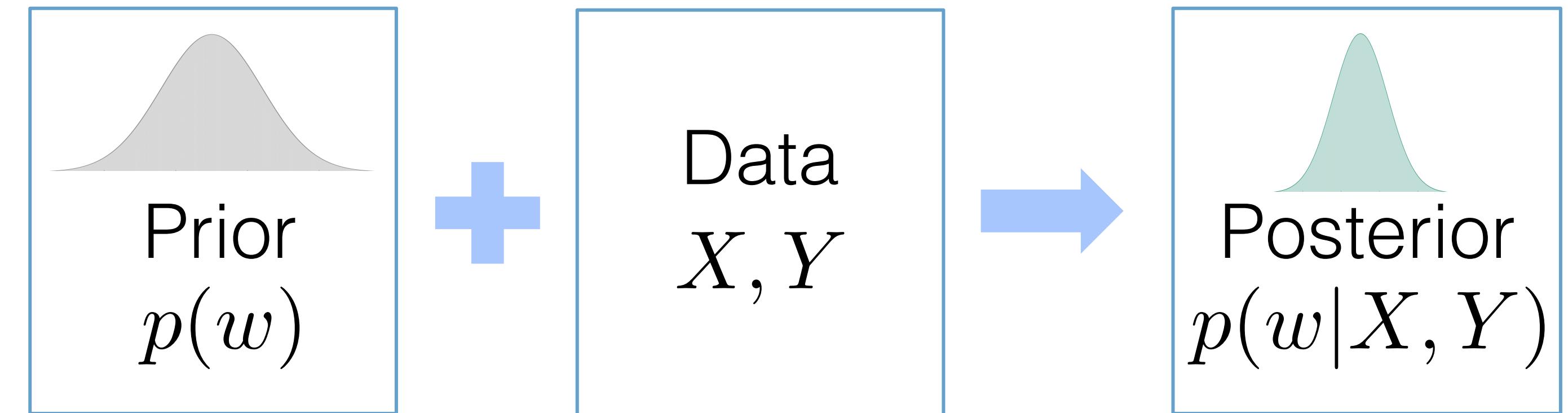
$$KL(q(w|\lambda) || p(w|X, Y)) \rightarrow \min_{\lambda}$$

Training Bayesian neural networks

Stochastic weights:



Bayesian Inference:



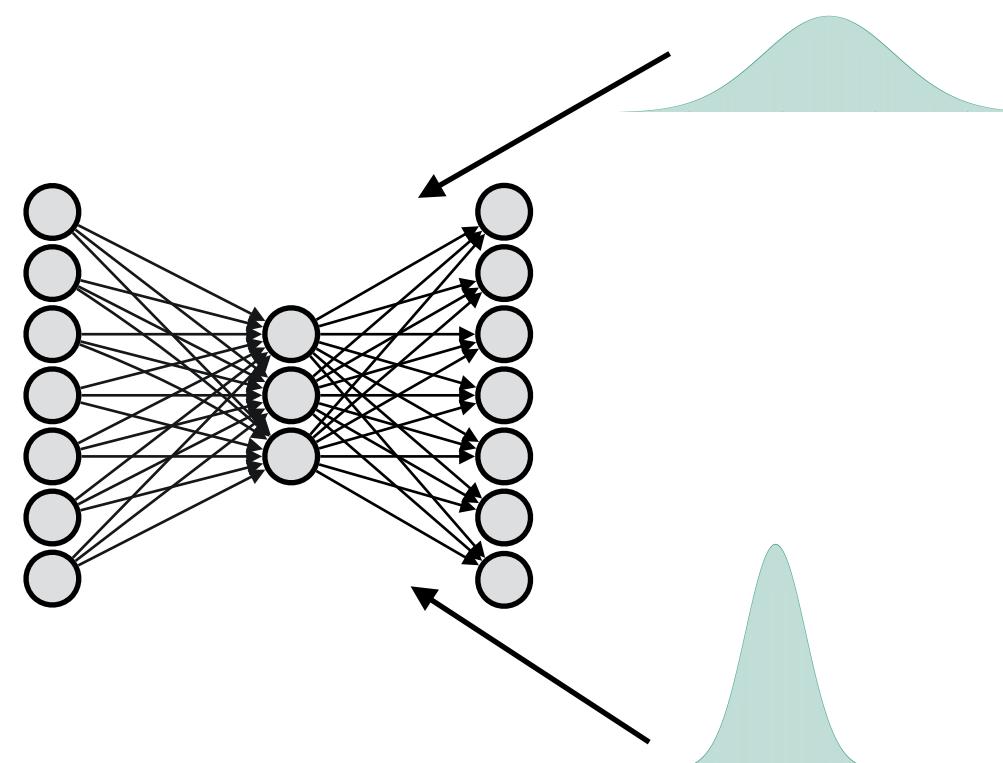
Posterior is intractable in neural networks → approximate it with $q(w|\lambda)$:

$$KL(q(w|\lambda) || p(w|X, Y)) \rightarrow \min_{\lambda}$$

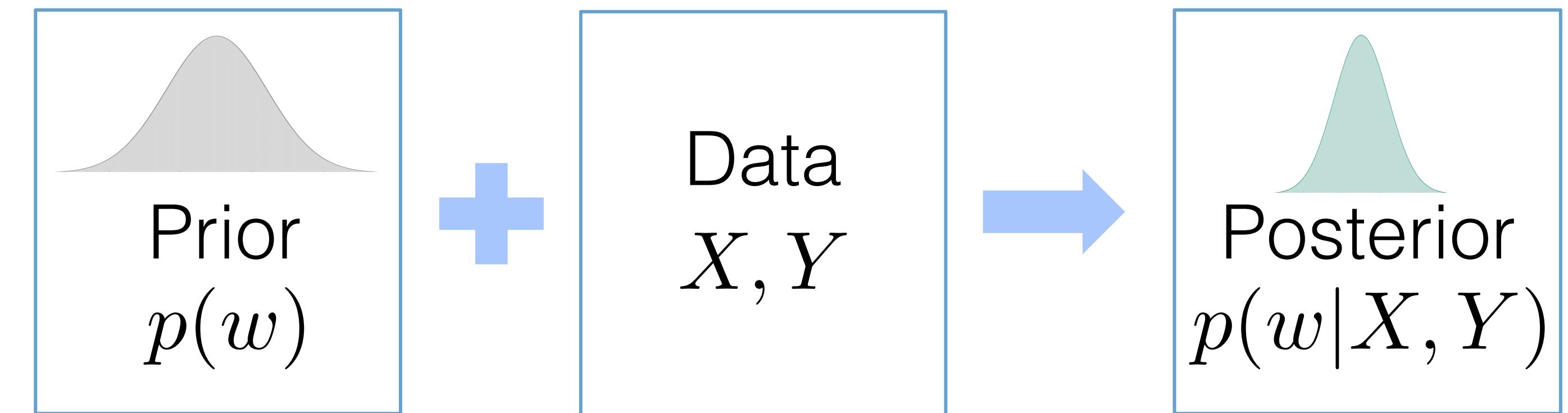
Intractable!

Training Bayesian neural networks

Stochastic weights:



Bayesian Inference:

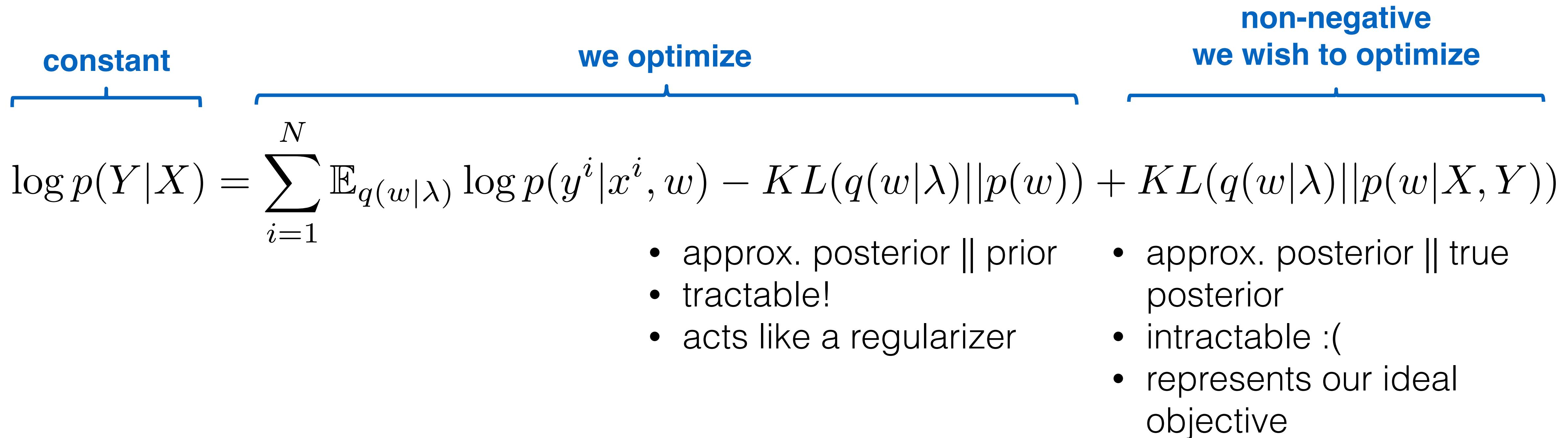


Equivalently, we can optimize ELBO to find approximate posterior $q(w|\lambda)$:

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w)}_{\text{Data term}} - \underbrace{KL(q(w|\lambda)||p(w))}_{\text{Regularizer}} \rightarrow \max_{\lambda}$$

KL-divergence: note there are two of them!

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w)}_{\text{Data term}} - \underbrace{KL(q(w|\lambda)||p(w))}_{\text{Regularizer}} \rightarrow \max_{\lambda}$$



Doubly stochastic variational inference

How to estimate the data term?

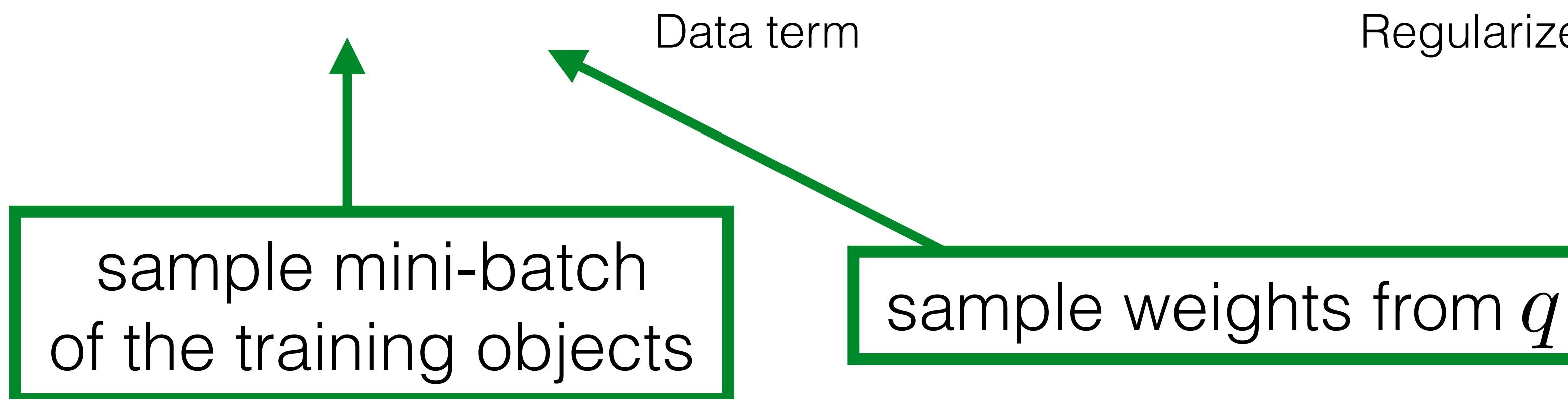
$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

Data term
??? Regularizer

Doubly stochastic variational inference

How to estimate the data term?

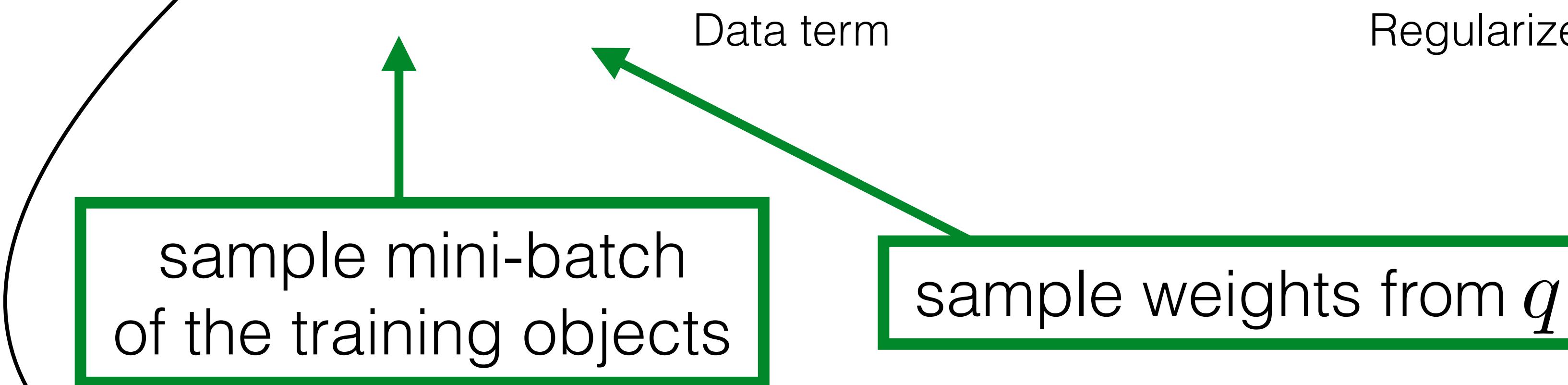
$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$



Doubly stochastic variational inference

How to estimate the data term?

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$



$$\sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, \hat{w}_j), \quad \hat{w}_j \sim q(w|\lambda) \quad i_j \sim \text{Unif}(1, \dots, N)$$

Doubly stochastic variational inference

How to estimate the data term?

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

sample mini-batch
of the training objects

sample weights from q

$$\sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, \hat{w}_j), \quad \hat{w}_j \sim q(w|\lambda) \quad i_j \sim \text{Unif}(1, \dots, N)$$

$\lambda?$ **How to differentiate?**

Doubly stochastic variational inference

How to estimate the data term?

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

$\frac{\partial}{\partial \lambda} ?$

sample mini-batch
of the training objects

sample weights from q

$$\sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, \hat{w}_j), \quad \hat{w}_j \sim q(w|\lambda) \quad i_j \sim \text{Unif}(1, \dots, N)$$

$\lambda? \quad \text{How to differentiate?}$

Gradient of expectation

$$\frac{\partial}{\partial \lambda} \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) = \frac{\partial}{\partial \lambda} \int q(w|\lambda) \log p(y^i|x^i, w) dw =$$

$$= \int \left(\frac{\partial}{\partial \lambda} q(w|\lambda) \right) \log p(y^i|x^i, w) dw$$

Not an expectation over q anymore!

Most popular solution: **reparametrization trick**

Reparametrization trick

As an example, consider 1-dim normal approximate posterior:

$$q(w|\lambda) = \mathcal{N}(\mu, \sigma^2), \quad \lambda = \{\mu, \sigma\}$$

Reparametrization trick

As an example, consider 1-dim normal approximate posterior:

$$q(w|\lambda) = \mathcal{N}(\mu, \sigma^2), \quad \lambda = \{\mu, \sigma\}$$

Property of normal distribution:

$$w \sim \mathcal{N}(\mu, \sigma^2) \iff w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

Reparametrization trick

As an example, consider 1-dim normal approximate posterior:

$$q(w|\lambda) = \mathcal{N}(\mu, \sigma^2), \quad \lambda = \{\mu, \sigma\}$$

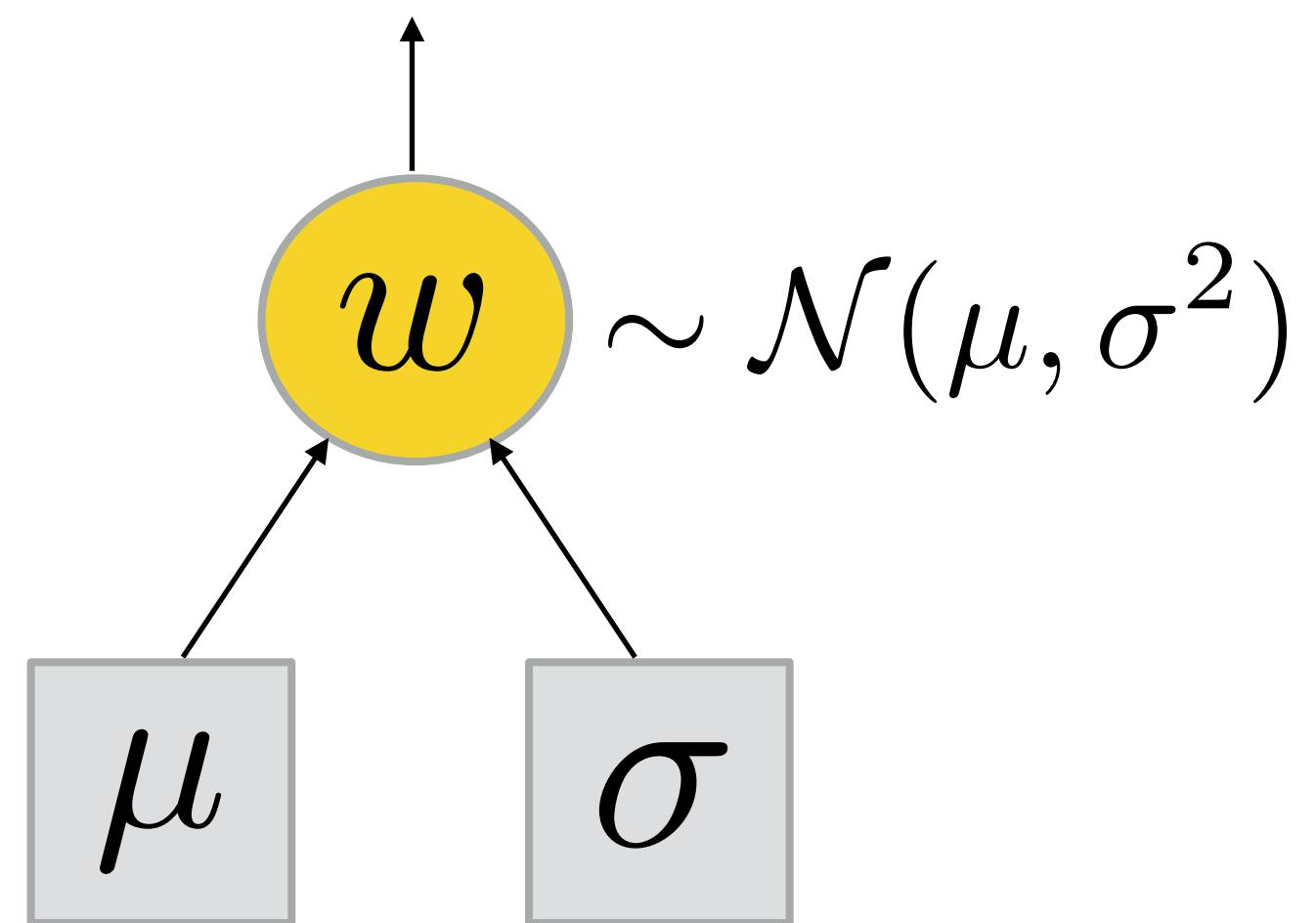
Property of normal distribution:

$$w \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

**distribution conditioned
on parameters** **unconditioned
distribution**

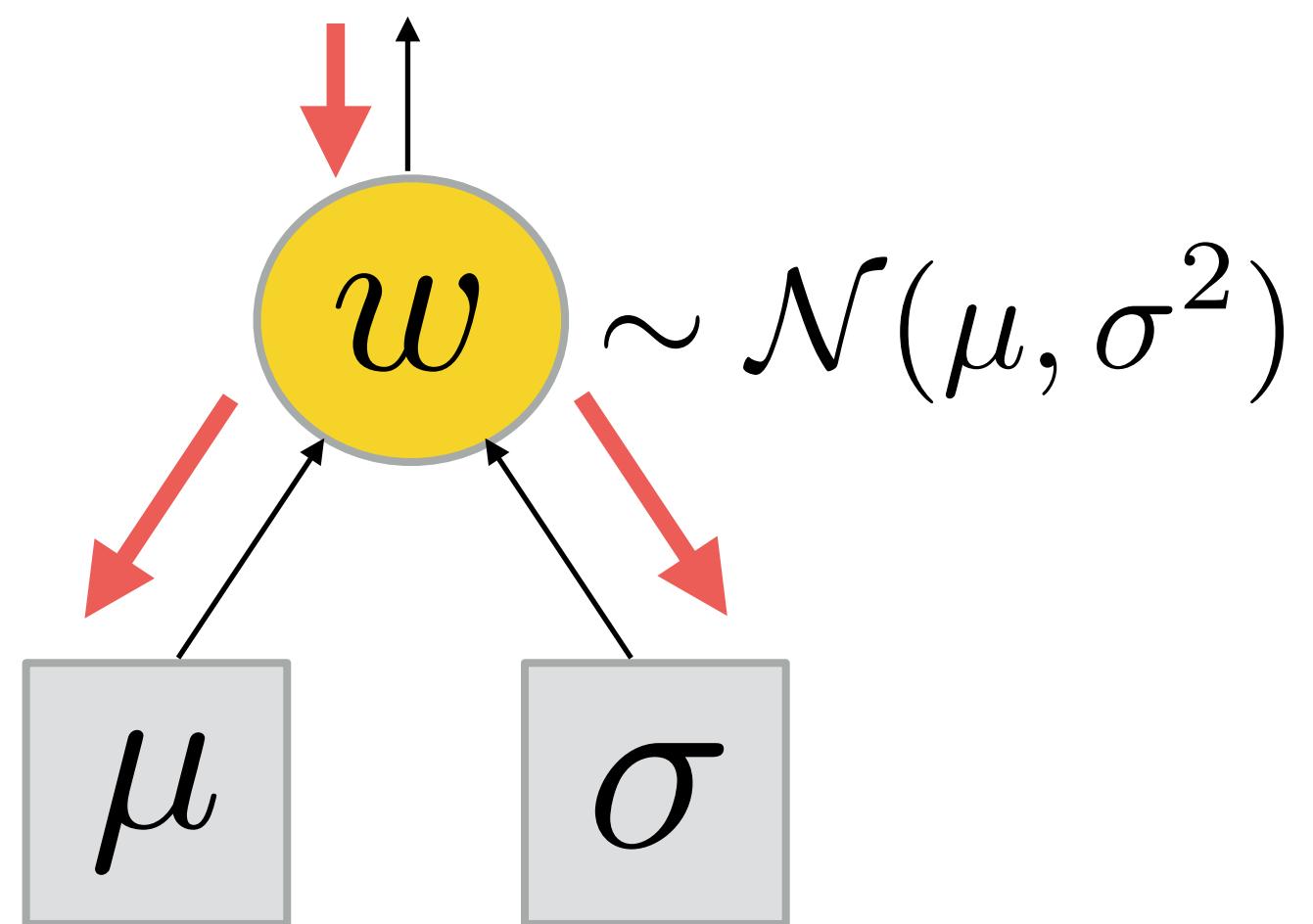
Reparametrization trick

$$w \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$



Reparametrization trick

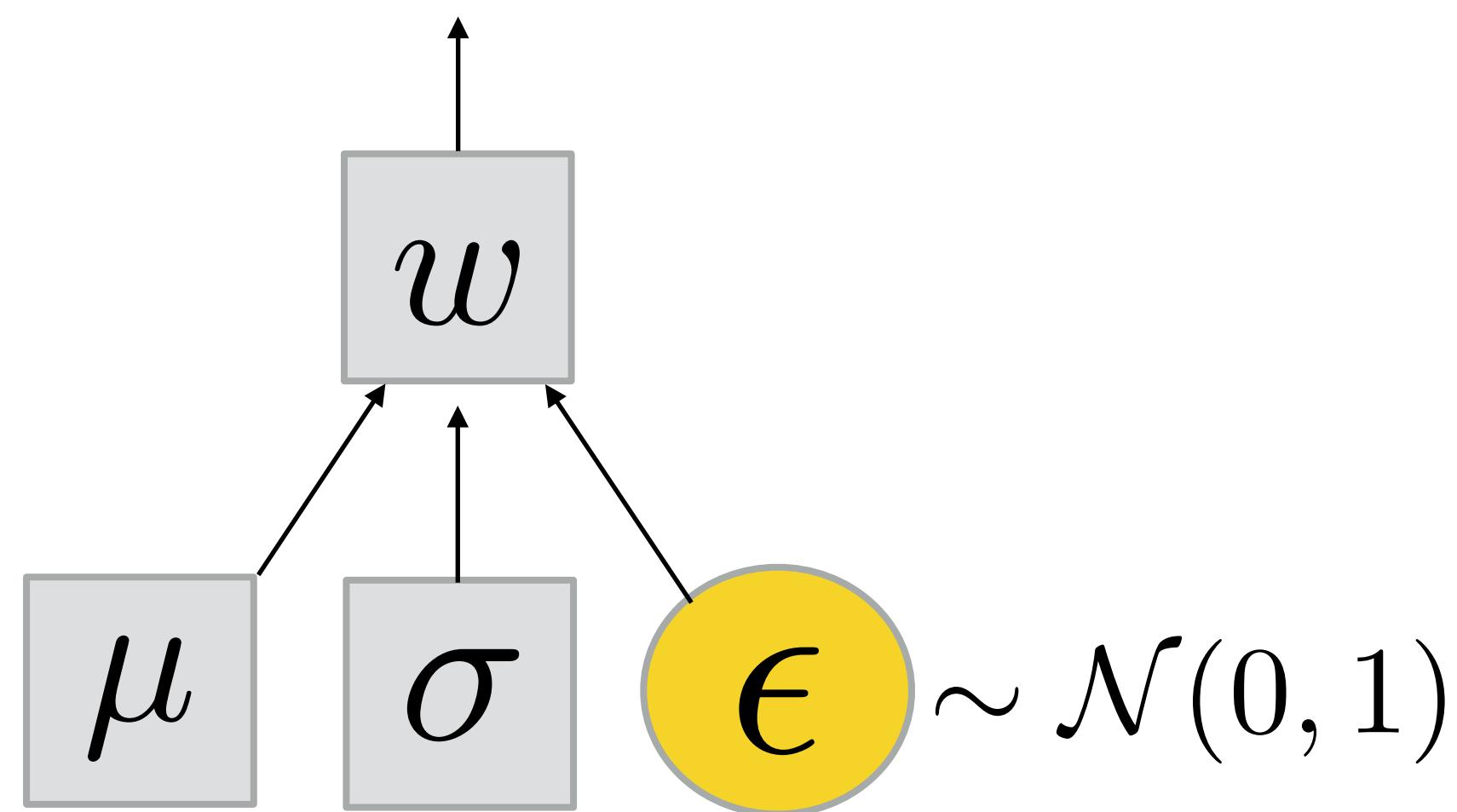
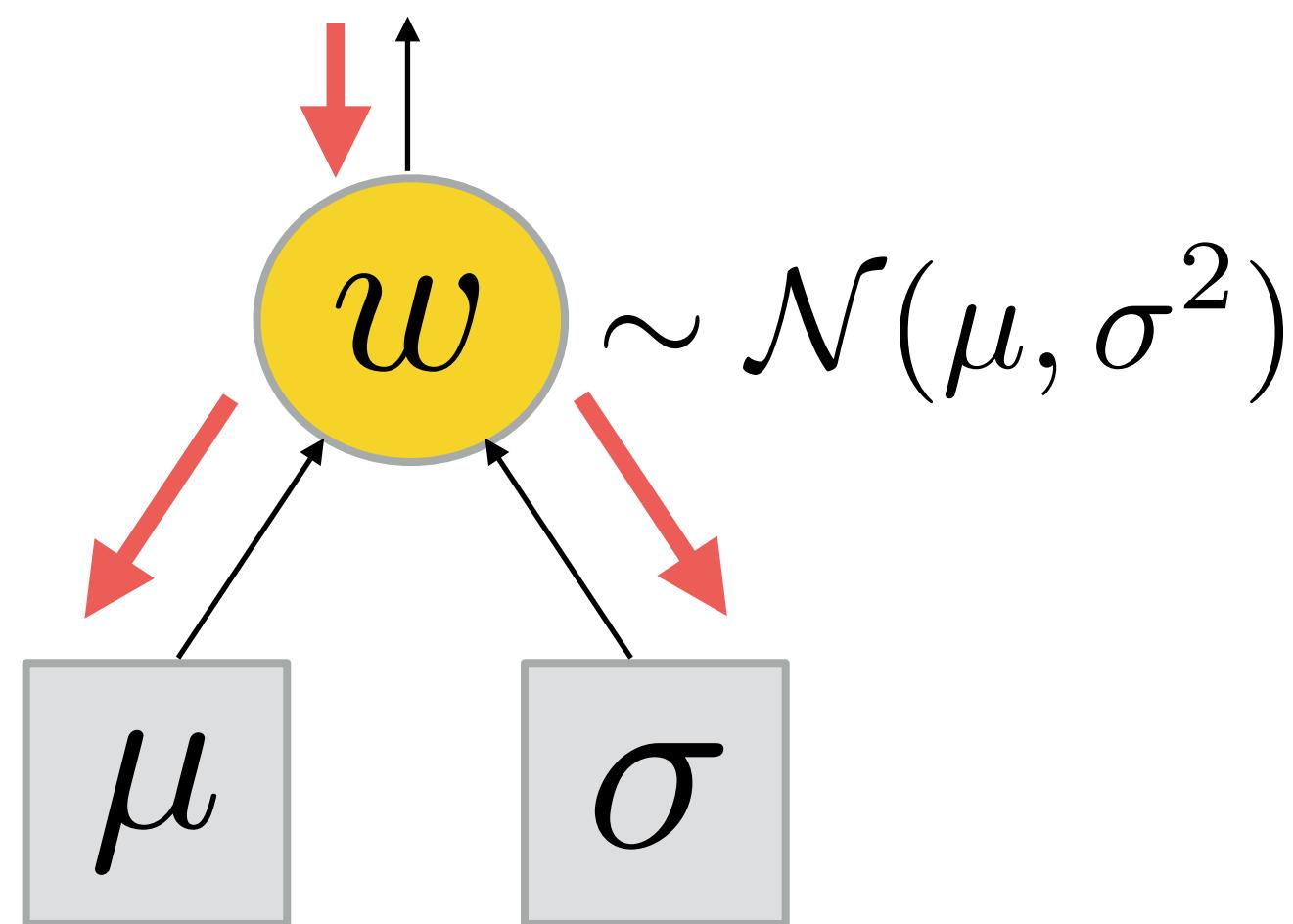
$$w \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$



**Gradients propagate
through randomness**

Reparametrization trick

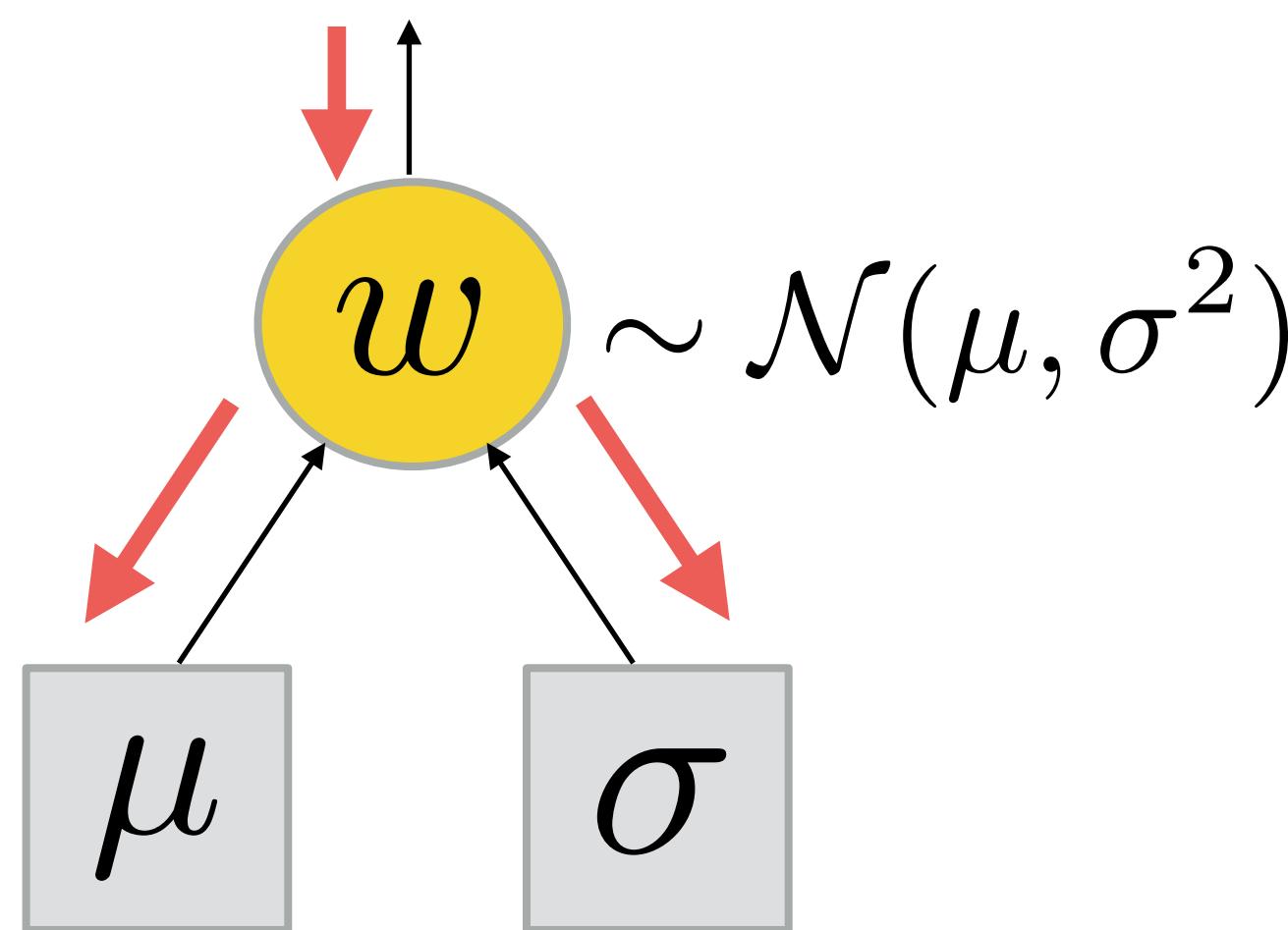
$$w \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$



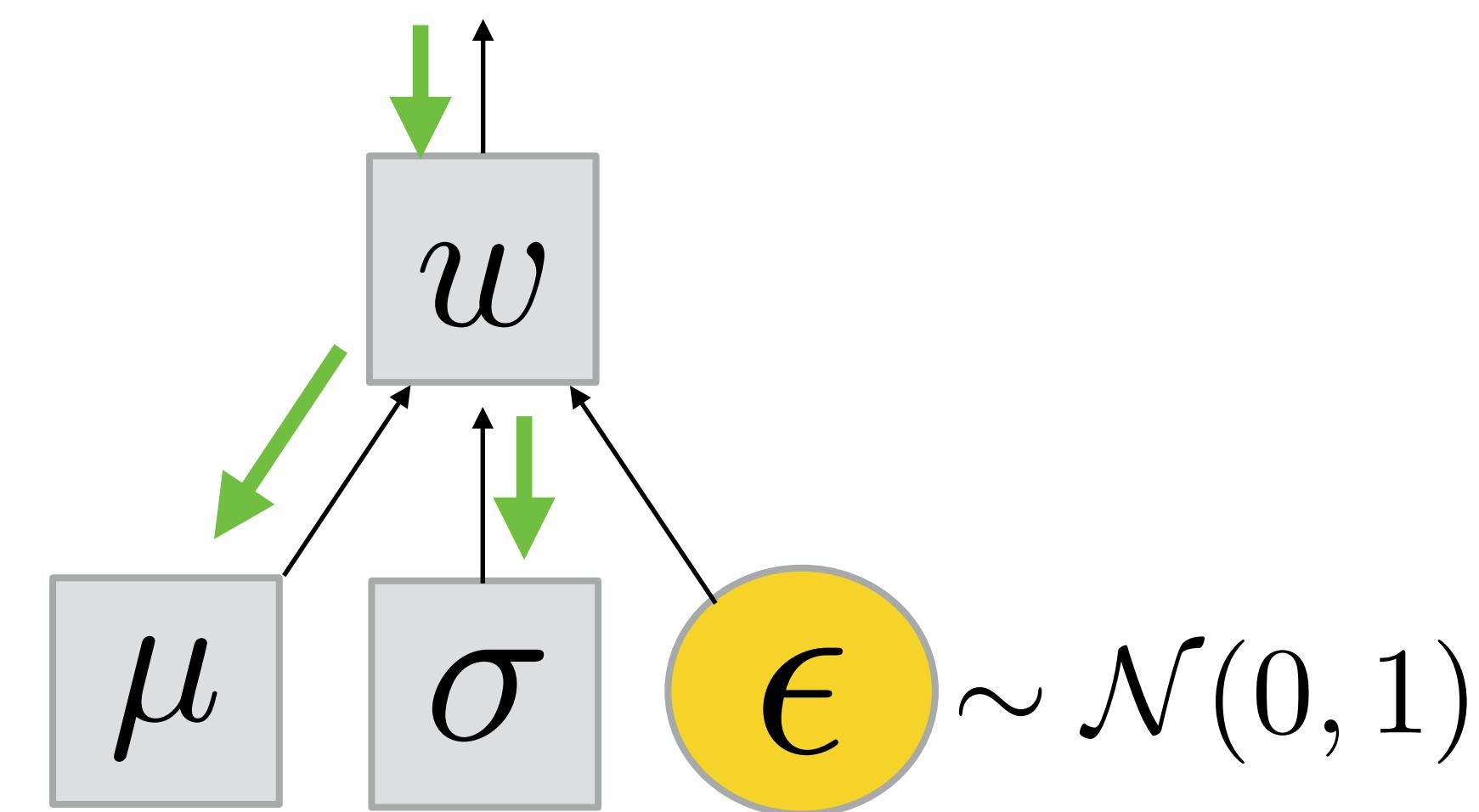
**Gradients propagate
through randomness**

Reparametrization trick

$$w \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$



**Gradients propagate
through randomness**



**Gradients propagate only
through deterministic nodes**

Reparametrization trick: general form

$$w \sim \mathcal{N}(\mu, \sigma^2) \Leftrightarrow w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

$$w \sim q(w|\lambda) \Leftrightarrow w = f(\lambda, \epsilon), \quad \epsilon \sim p(\epsilon)$$

**distribution
conditioned
on parameters**

**dependency
on parameters**

**unconditioned
distribution**

Reparametrization trick

$$w \sim q(w|\lambda) \Leftrightarrow w = f(\lambda, \epsilon), \quad \epsilon \sim p(\epsilon)$$

distribution conditioned on parameters **dependency on parameters** **unconditioned distribution**

$$\begin{aligned} \nabla_{\lambda} \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) &= \nabla_{\lambda} \mathbb{E}_{p(\epsilon)} \log p(y^i|x^i, w = f(\lambda, \epsilon)) = \\ &= \mathbb{E}_{p(\epsilon)} \nabla_{\lambda} \log p(y^i|x^i, w = f(\lambda, \epsilon)) \end{aligned}$$

Reparametrization trick: examples

$$q(w|\lambda) \longrightarrow w = f(\lambda, \epsilon), \quad p(\epsilon)$$

$q(w \lambda)$	$p(\epsilon)$	$f(\lambda, \epsilon)$
$\mathcal{N}(\mu, \sigma^2)$	$\mathcal{N}(0, 1)$	$w = \sigma\epsilon + \mu$
$\mathcal{G}(1, \beta)$	$\mathcal{G}(1, 1)$	$w = \beta\epsilon$
$\mathcal{E}(\lambda)$	$\mathcal{U}(0, 1)$	$w = -\frac{\log \epsilon}{\lambda}$
$\mathcal{N}(\mu, \Sigma)$	$\mathcal{N}(0, I)$	$w = A\epsilon + \mu, \text{ where } AA^T = \Sigma$

Reparametrization trick: examples

$$q(w|\lambda) \rightarrow w = f(\lambda, \epsilon), p(\epsilon)$$

$q(w \lambda)$	$p(\epsilon)$	$f(\lambda, \epsilon)$
$\mathcal{N}(\mu, \sigma^2)$	$\mathcal{N}(0, 1)$	$w = \sigma\epsilon + \mu$
$\mathcal{G}(1, \beta)$	$\mathcal{G}(1, 1)$	$w = \beta\epsilon$
$\mathcal{E}(\lambda)$	$\mathcal{U}(0, 1)$	$w = -\frac{\log \epsilon}{\lambda}$
$\mathcal{N}(\mu, \Sigma)$	$\mathcal{N}(0, I)$	$w = A\epsilon + \mu, \text{ where } AA^T = \Sigma$

noise on
weights!

Local reparametrization trick

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$w_{ij} = \mu_{ij} + \epsilon_{ij}\sigma_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(0, 1)$$

- Too expensive! ($|w| \times$ mini-batch size)
- One sample per mini-batch? High variance of stochastic gradients & correlated predictions

Local reparametrization trick

1. Reparametrization trick \Rightarrow unbiased estimate of the gradients

$$w_{ij} = \mu_{ij} + \epsilon_{ij}\sigma_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(0, 1)$$

2. Local reparametrization trick: **sample preactivations** instead of weights
 \Rightarrow reduced variance of the gradients & uncorrelated predictions

$$w_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2) \quad \mathbb{E}X_2 = MX_1 \quad \text{Var}X_2 = \Sigma^2 X_1^2$$

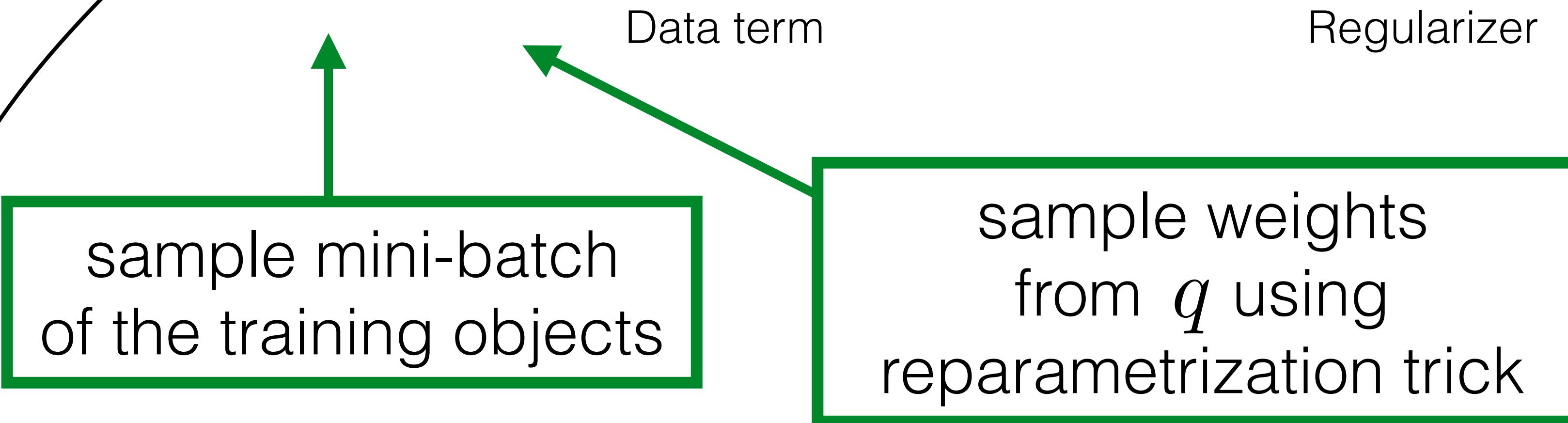
$$X_2 = WX_1 \quad X_2 \sim \mathcal{N}(MX_1, \Sigma^2 X_1^2) \quad \text{blue means element-wise}$$

$$X_2 = MX_1 + \sqrt{\Sigma^2 X_1^2} \odot \epsilon$$

Training: putting everything together

How to estimate the data term?

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$



$$\sum_{j=1}^m \log p(y^{i_j}|x^{i_j}, w = f(\lambda, \epsilon^j)), \quad \epsilon^j \sim p(\epsilon) \quad i_j \sim \text{Unif}(1, \dots, N)$$

Training: putting everything together

Deterministic neural network:

$$w^{new} = w^{old} + \eta \frac{\partial}{\partial w} \sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, w^{old}) \quad i_j \sim \text{Unif}(1, \dots, N)$$

Bayesian neural network:

$$\lambda^{new} = \lambda^{old} + \eta \frac{\partial}{\partial \lambda} \sum_{j=1}^m \log p(y^{i_j} | x^{i_j}, w = f(\lambda^{old}, \epsilon^j)), \quad i_j \sim \text{Unif}(1, \dots, N) \\ \epsilon^j \sim p(\epsilon)$$

m — mini-batch size

η — learning rate

Questions

Questions

- Assume $|w| = M$. If we approximate the posterior with a fully-factorized normal distribution, how many parameters do we train?

Questions

- Assume $|w| = M$. If we approximate the posterior with a fully-factorized normal distribution, how many parameters do we train?
- If we approximate the posterior with an arbitrary normal distribution (full covariance matrix)?

Questions

- Assume $|w| = M$. If we approximate the posterior with a fully-factorized normal distribution, how many parameters do we train?
- If we approximate the posterior with an arbitrary normal distribution (full covariance matrix)?
- What can we do if we cannot compute the KL-regularizer?

$$KL(q(w|\lambda) || p(w)) = \int q(w|\lambda) \log \frac{q(w|\lambda)}{p(w)} dw$$

Questions

- Assume $|w| = M$. If we approximate the posterior with a fully-factorized normal distribution, how many parameters do we train?
- If we approximate the posterior with an arbitrary normal distribution (full covariance matrix)?
- What can we do if we cannot compute the KL-regularizer?

$$KL(q(w|\lambda) || p(w)) = \int q(w|\lambda) \log \frac{q(w|\lambda)}{p(w)} dw \approx \log \frac{q(\hat{w}|\lambda)}{p(\hat{w})}, \quad \hat{w} \sim q(w|\lambda)$$

+ reparametrization trick

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

Model specification:

- Choose particular prior

Training:

- Choose particular family for approximate posterior
- How to compute the KL-divergence?

Software

- Pyro (based on PyTorch)
- TensorFlow Probability (based on TensorFlow)
- Edward (based on TensorFlow)
- PyMC (based on Theano, pre-release with TensorFlow Probability)
- <https://github.com/JavierAntoran/Bayesian-Neural-Networks> —
PyTorch implementations of popular Bayesian deep learning papers

Summary

- A lot of BNN advantages: regularization, ensembling, uncertainty estimation, ...
- To train BNN, one should optimize ELBO using DSVI & RT
- Three steps towards a particular method

Plan

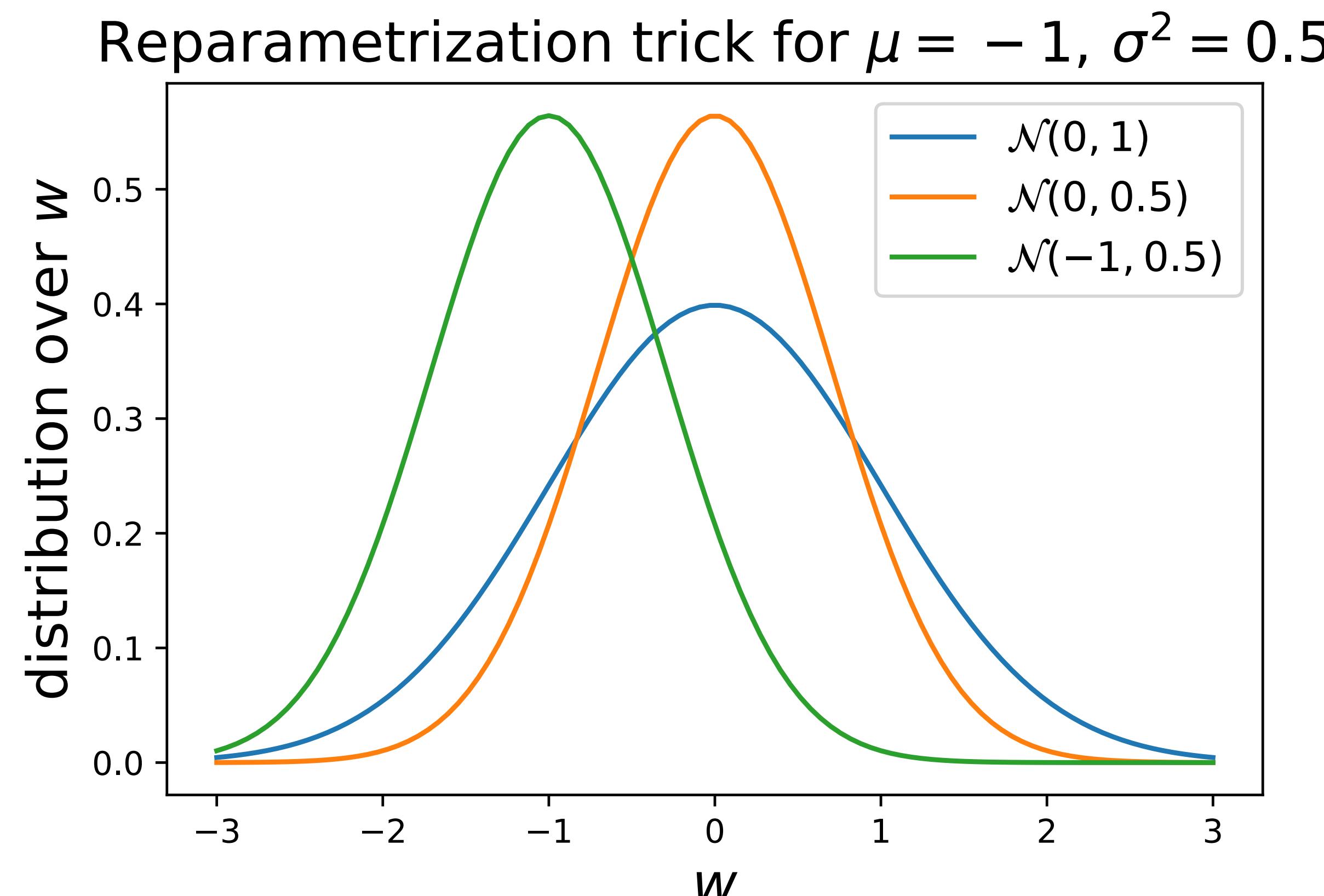
- Advantages of using Bayesian neural networks
- Training Bayesian neural networks
- Q&A + exercises

Questions

Reparametrization trick: recap

$$w \sim \mathcal{N}(\mu, \sigma^2) \quad \Leftrightarrow \quad w = \mu + \sigma\epsilon, \quad \epsilon \sim \mathcal{N}(0, 1)$$

$$w \sim q(w|\lambda) \quad \Leftrightarrow \quad w = f(\lambda, \epsilon), \quad \epsilon \sim p(\epsilon)$$

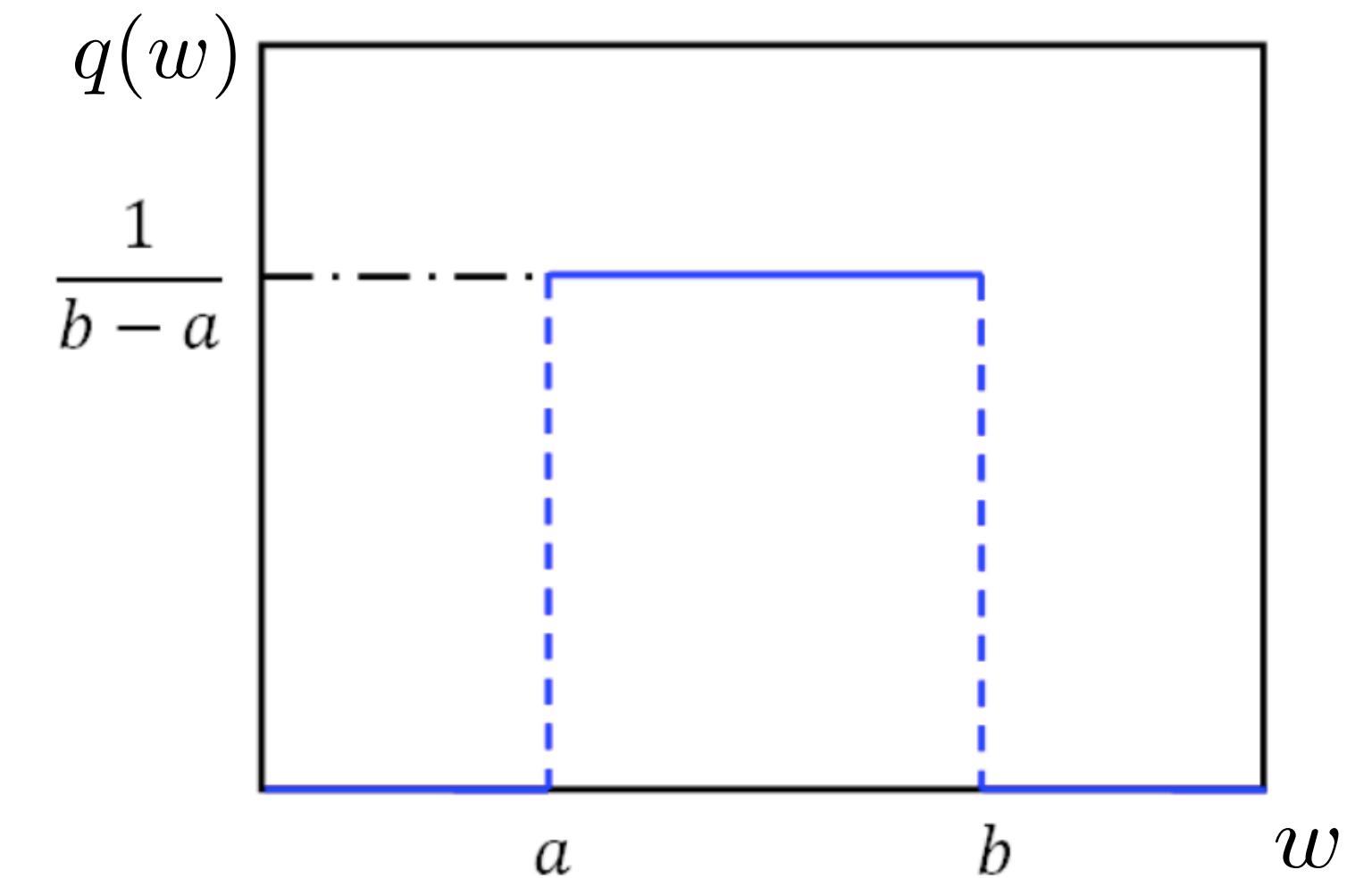


Exercises

Reparametrization trick: $w \sim q(w|\lambda) \Leftrightarrow w = f(\lambda, \epsilon), \epsilon \sim p(\epsilon)$

- How can we reparametrize uniform distribution?

$$q(w|a, b) = \begin{cases} \frac{1}{b-a}, & w \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$



- How can we reparametrize a mixture of normal distributions?

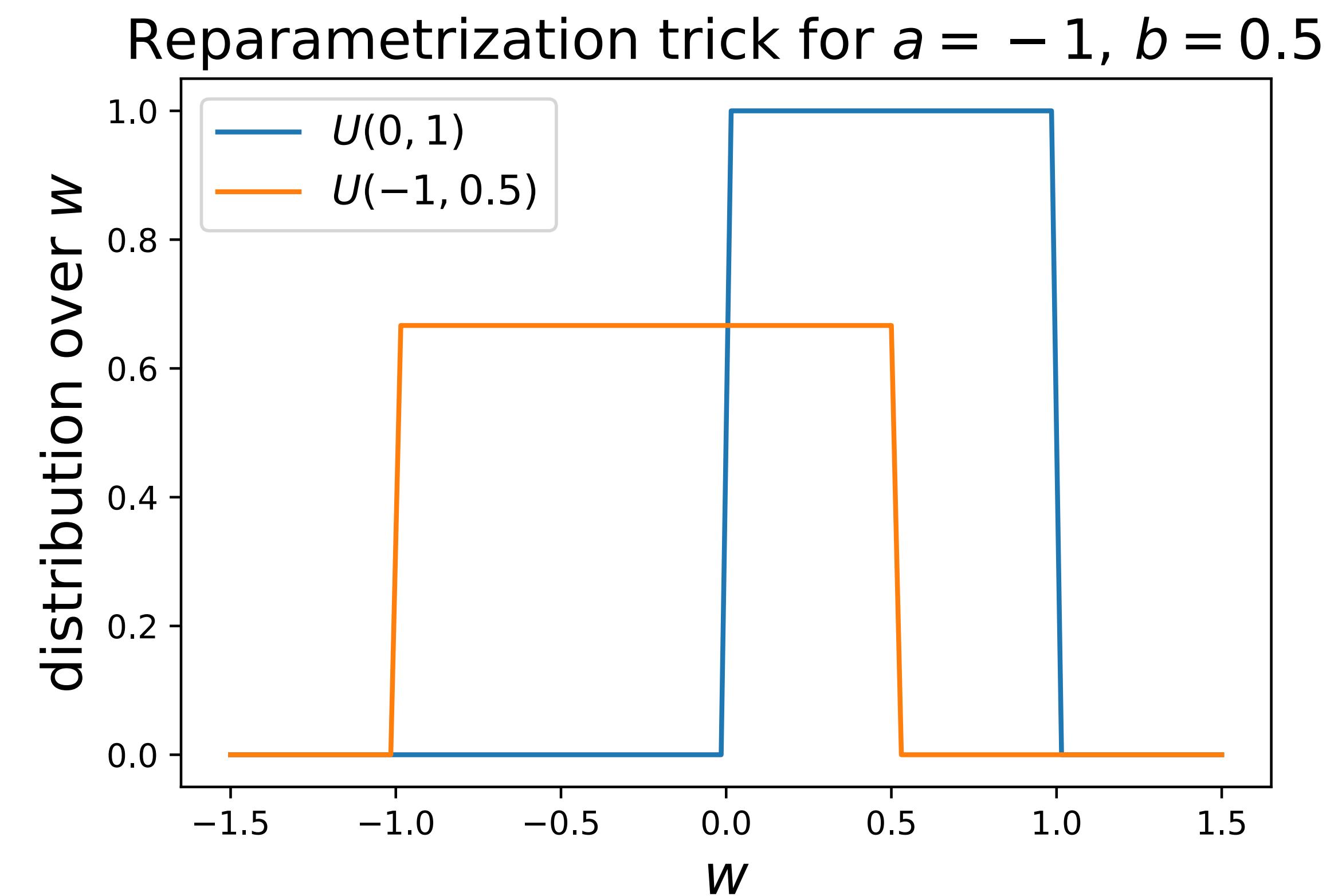
$$q(w|\mu, \sigma) = \sum_{s=1}^S \pi_s \mathcal{N}(w|\mu_s, \sigma_s), \quad \sum_{s=1}^S \pi_s = 1, \pi_s \geq 0 \quad w \in \mathbb{R}$$

(π is fixed)

Reparametrizing uniform distribution

$$q(w|a, b) = \begin{cases} \frac{1}{b-a}, & w \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$

$w \sim q(w|a, b) \Leftrightarrow$
 $w = (1 - t)a + tb, t \sim U(0, 1)$



Reparametrizing a mixture of normal distributions

$$q(w|\mu, \sigma) = \sum_{s=1}^S \pi_s \mathcal{N}(w|\mu_s, \sigma_s), \quad \sum_{s=1}^S \pi_s = 1, \pi_s \geq 0$$

$$\begin{aligned} w \sim q(w|\mu, \sigma) &\iff w = f(\mu, \sigma, \epsilon, z) = \mu_z + \epsilon \sigma_z \\ \epsilon \sim \mathcal{N}(\epsilon|0, 1), \quad z &\sim Cat(\pi_1, \dots, \pi_S) \end{aligned}$$

Example: binary dropout

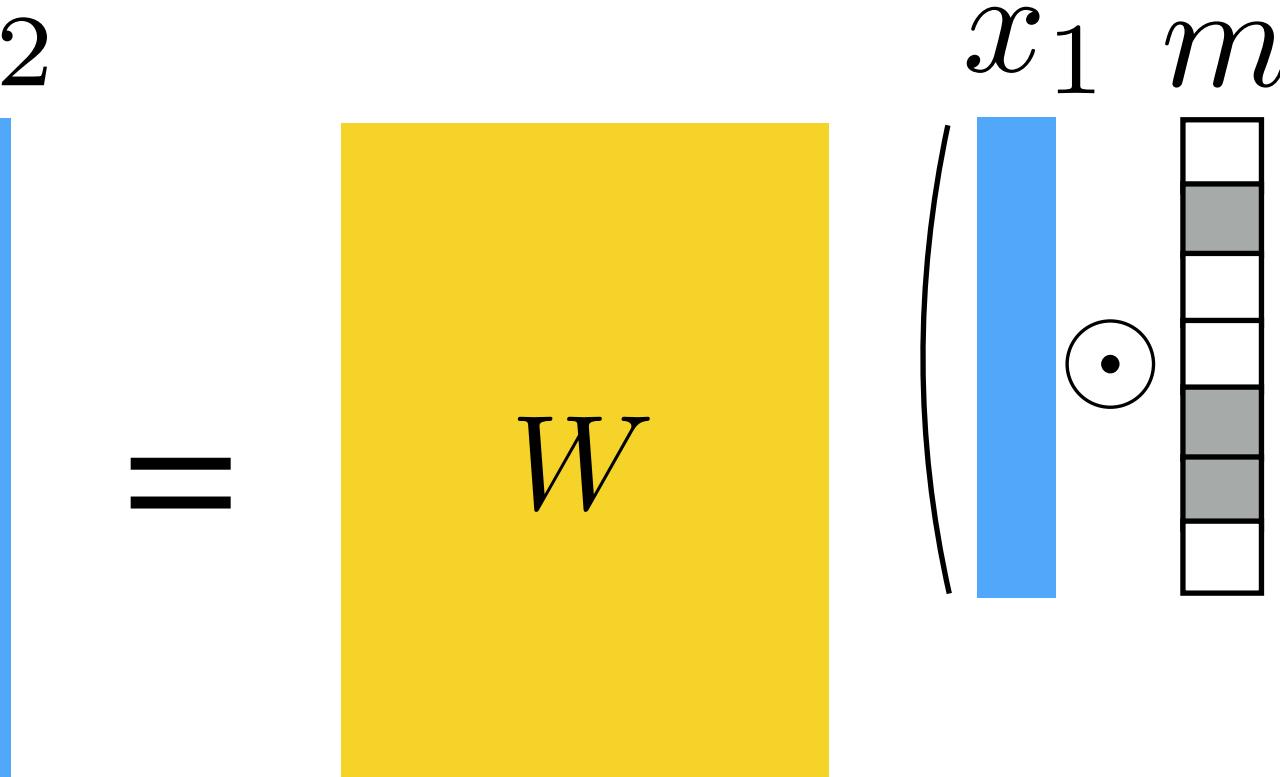
Linear layer: $x_2 = Wx_1$, W — weight matrix

With dropout: $x_2 = W(x_1 \odot m)$, $m_i \sim \text{Bernoulli}(p)$, p — dropout rate

Example: binary dropout

Linear layer: $x_2 = Wx_1$, W — weight matrix

With dropout: $x_2 = W(x_1 \odot m)$, $m_i \sim \text{Bernoulli}(p)$, p — dropout rate

$$x_2 = W \begin{pmatrix} x_1 & m \\ \odot & \end{pmatrix}$$


Example: binary dropout

Linear layer: $x_2 = Wx_1$, W — weight matrix

With dropout: $x_2 = W(x_1 \odot m)$, $m_i \sim \text{Bernoulli}(p)$, p — dropout rate

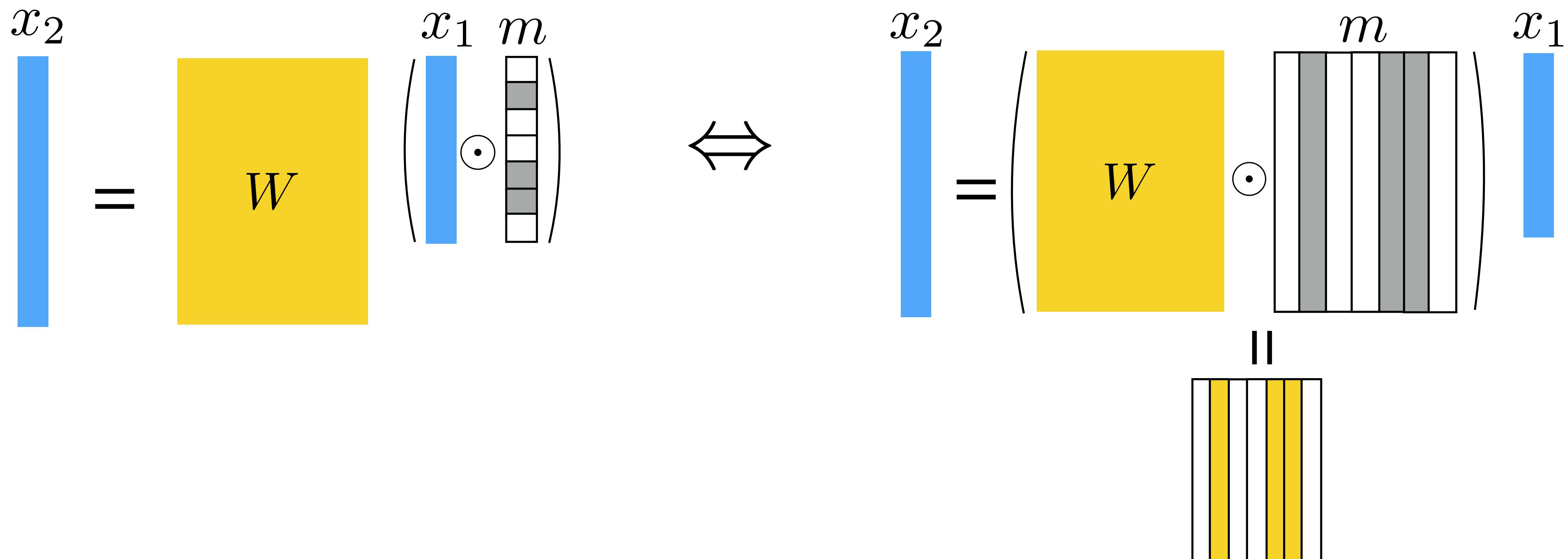
$$x_2 = \begin{matrix} \text{blue bar} \\ = \end{matrix} \quad \boxed{W} \quad \left(\begin{matrix} x_1 & m \\ \odot \end{matrix} \right) \quad \iff \quad \begin{matrix} x_2 \\ = \end{matrix} \quad \boxed{W} \quad \odot \quad \boxed{\begin{matrix} m \\ \odot \\ x_1 \end{matrix}}$$

The diagram illustrates the equivalence between a standard linear layer and one with binary dropout. On the left, a blue bar labeled x_2 is followed by an equals sign, then a yellow box labeled W , and finally a matrix multiplication symbol. This is followed by a circled dot symbol (\odot) and a matrix with columns labeled x_1 and m . An equivalence arrow (\iff) points to the right side of the diagram. On the right, a blue bar labeled x_2 is followed by an equals sign, then a yellow box labeled W , and finally a circled dot symbol (\odot). This is followed by a large matrix with columns labeled m and x_1 . The matrix has vertical gray bars in the m column and white bars in the x_1 column, representing binary dropout.

Example: binary dropout

Linear layer: $x_2 = Wx_1$, W — weight matrix

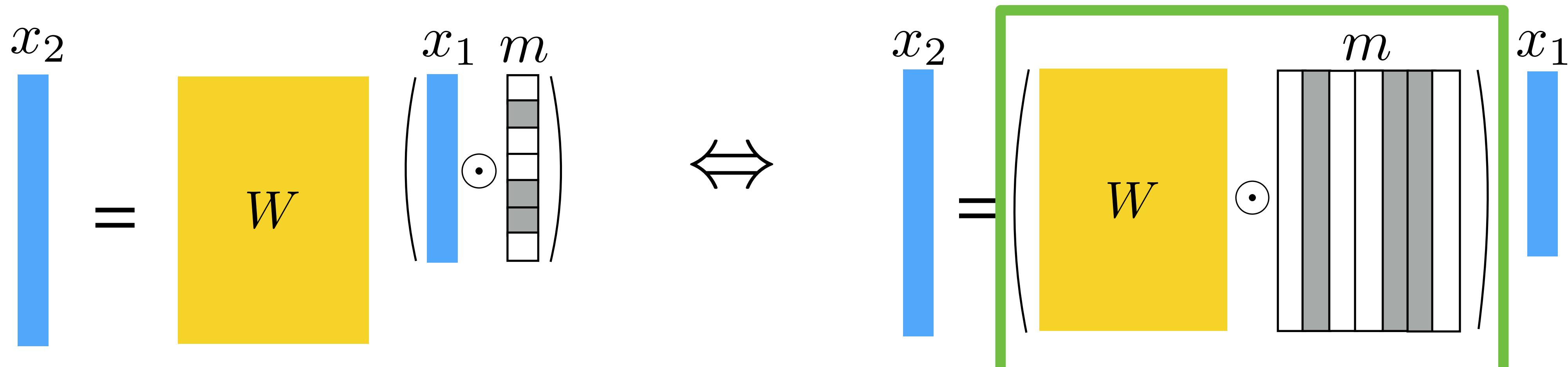
With dropout: $x_2 = W(x_1 \odot m)$, $m_i \sim \text{Bernoulli}(p)$, p — dropout rate



Example: binary dropout

Linear layer: $x_2 = Wx_1$, W — weight matrix

With dropout: $x_2 = W(x_1 \odot m)$, $m_i \sim \text{Bernoulli}(p)$, p — dropout rate

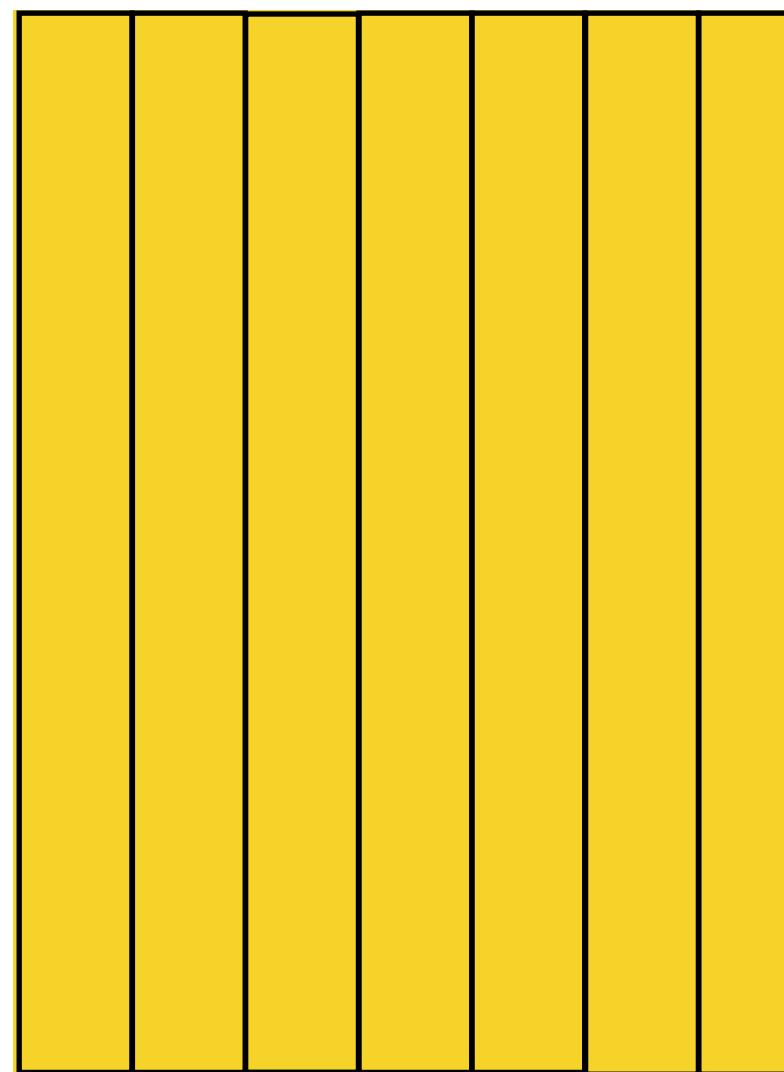


Applying dropout means
sampling weights!

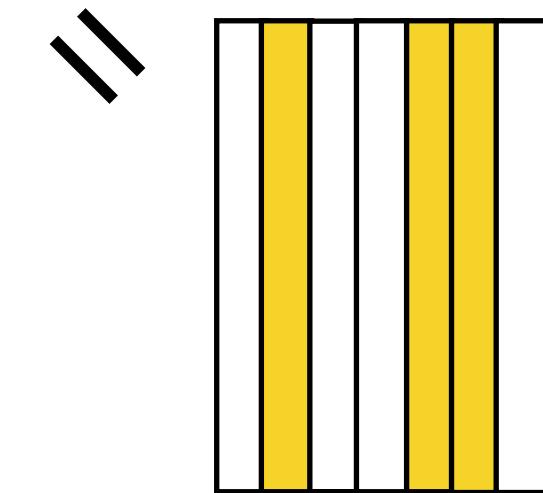
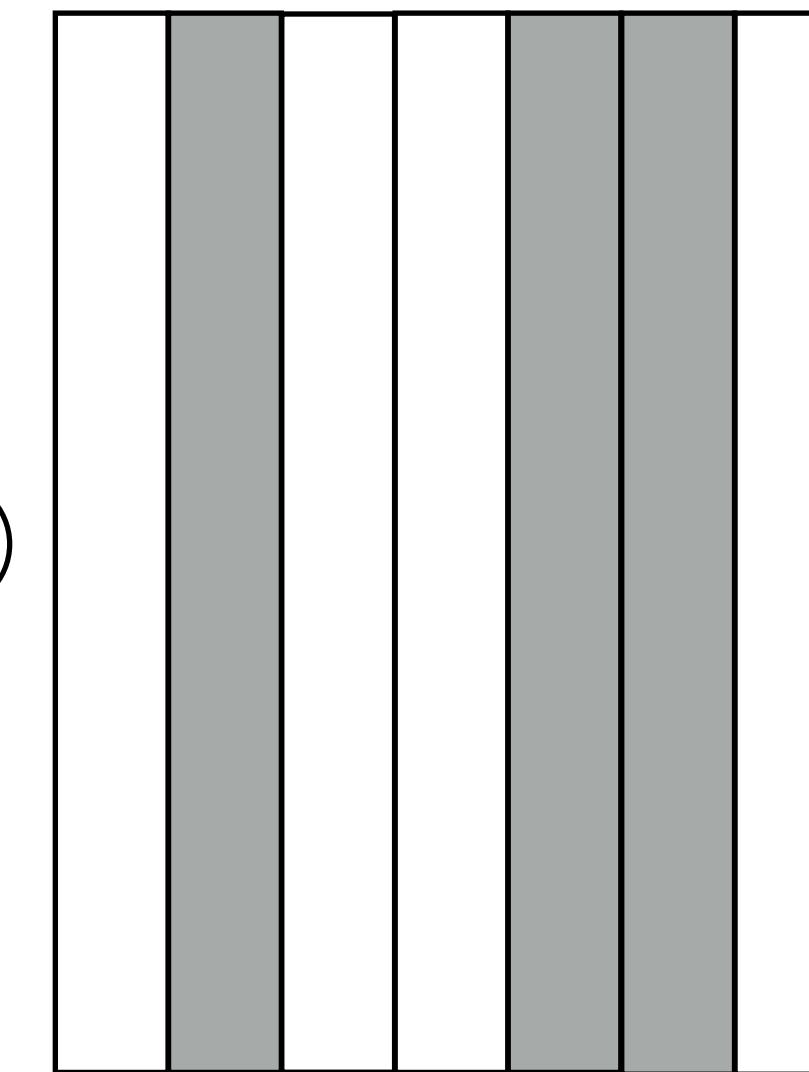
Example: binary dropout

Weight matrix:

$\mu_1 \mu_2 \mu_3 \dots$



$m_1 m_2 m_3 \dots$



reparametrization

$$w_i = f(\mu_i, m_i) = \mu_i \cdot m_i$$
$$m_i \sim \text{Bernoulli}(p)$$

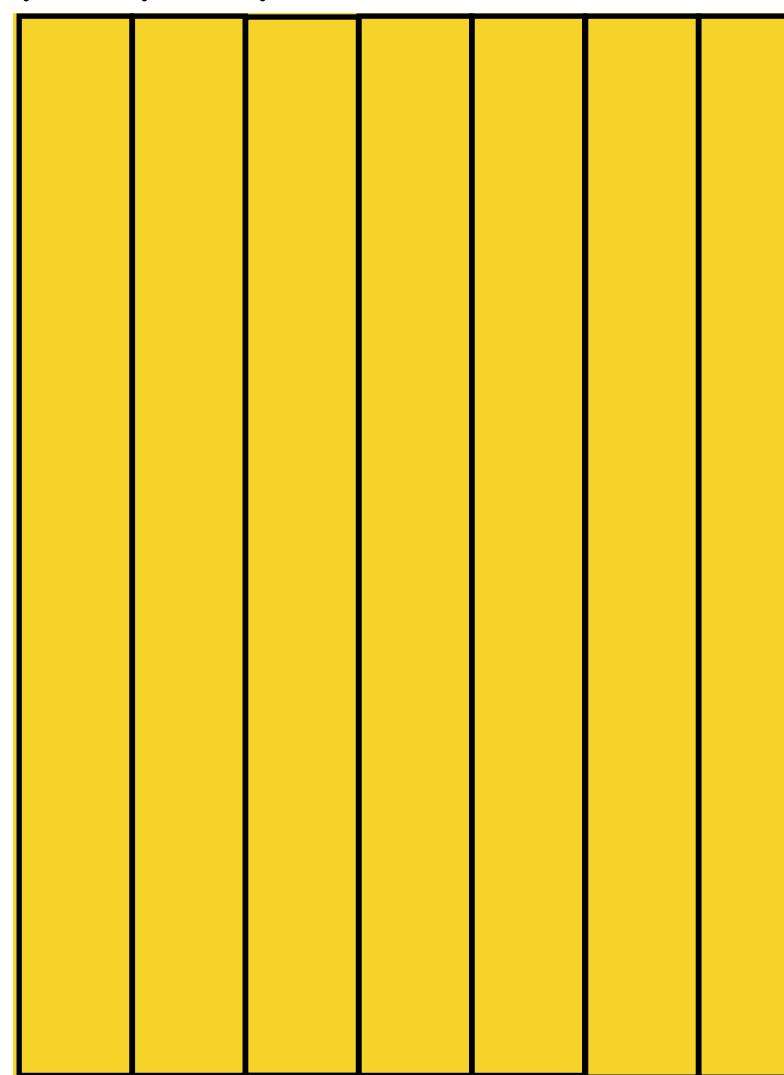
← 0 or 1!

p is a fixed hyperparameter!

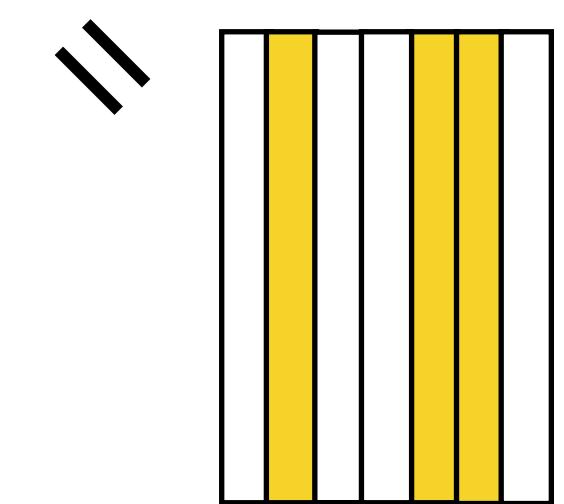
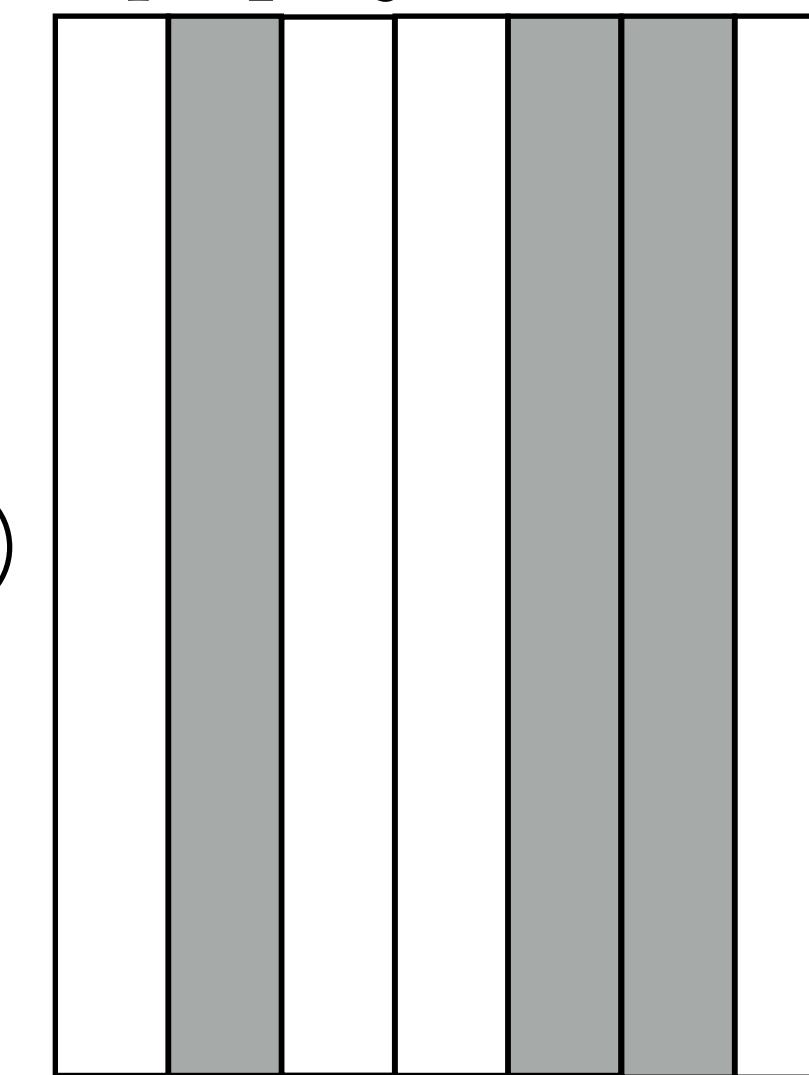
Example: binary dropout

Weight matrix:

$\mu_1 \mu_2 \mu_3 \dots$



$m_1 m_2 m_3 \dots$



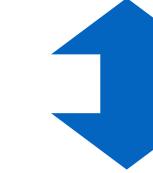
$$w_i = f(\mu_i, m_i) = \mu_i \cdot m_i$$

$$m_i \sim \text{Bernoulli}(p)$$

reparametrization

$$q(W|\mu) = \prod_i q(w_i|\mu_i)$$

$$q(w_i|\mu_i) = p \delta(0) + (1 - p) \delta(\mu_i)$$



p is a fixed hyperparameter!

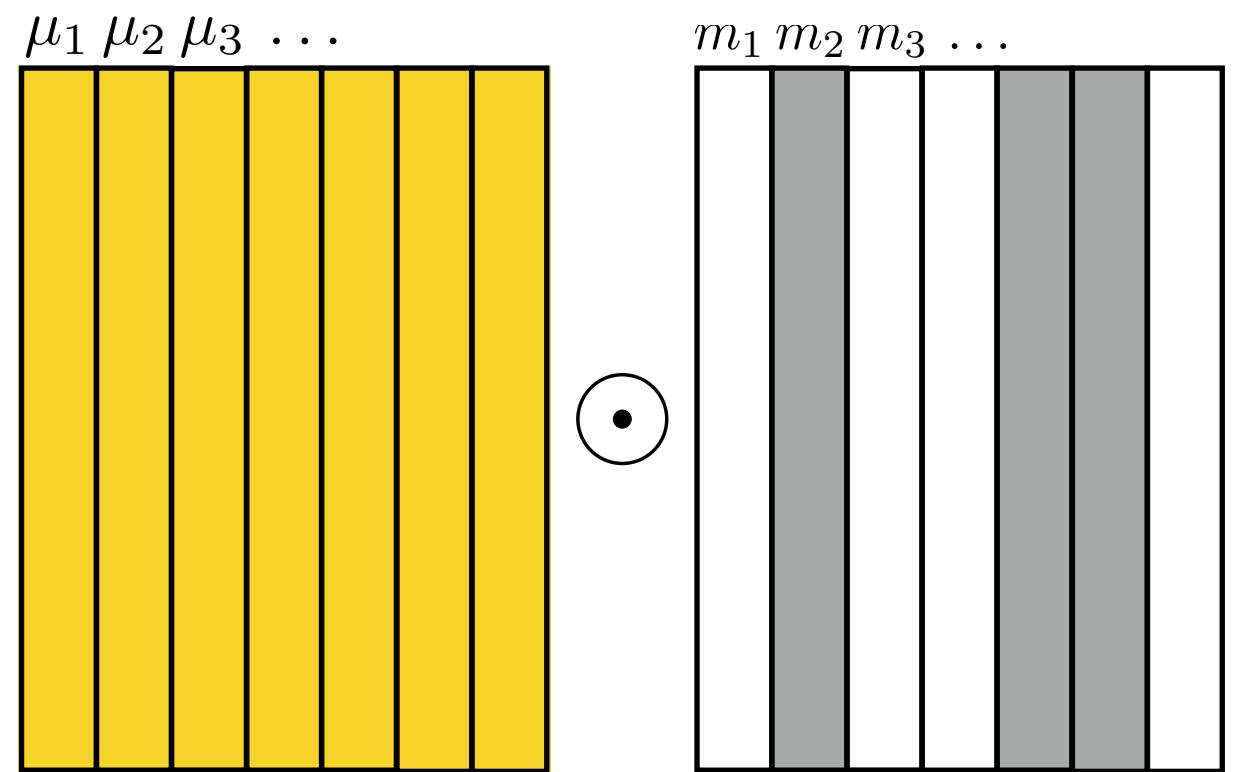
Example: binary dropout

Prior: ?

Approximate posterior: $q(W|\mu) = \prod_i q(w_i|\mu_i)$

$$q(w_i|\mu_i) = p \delta(0) + (1 - p) \delta(\mu_i)$$

Weight matrix:



Approximate KL-divergence: ?

Example: binary dropout

Prior: $p(W) = \prod_{i,j} p(w_{ij}), \quad p(w_{ij}) = \mathcal{N}(0, 1)$

Approximate posterior: $q(W|\mu) = \prod_i q(w_i|\mu_i)$

$$q(w_i|\mu_i) = p \delta(0) + (1 - p) \delta(\mu_i)$$

Approximate KL-divergence: ?

Example: binary dropout

Prior: $p(W) = \prod_{i,j} p(w_{ij}), \quad p(w_{ij}) = \mathcal{N}(0, 1)$

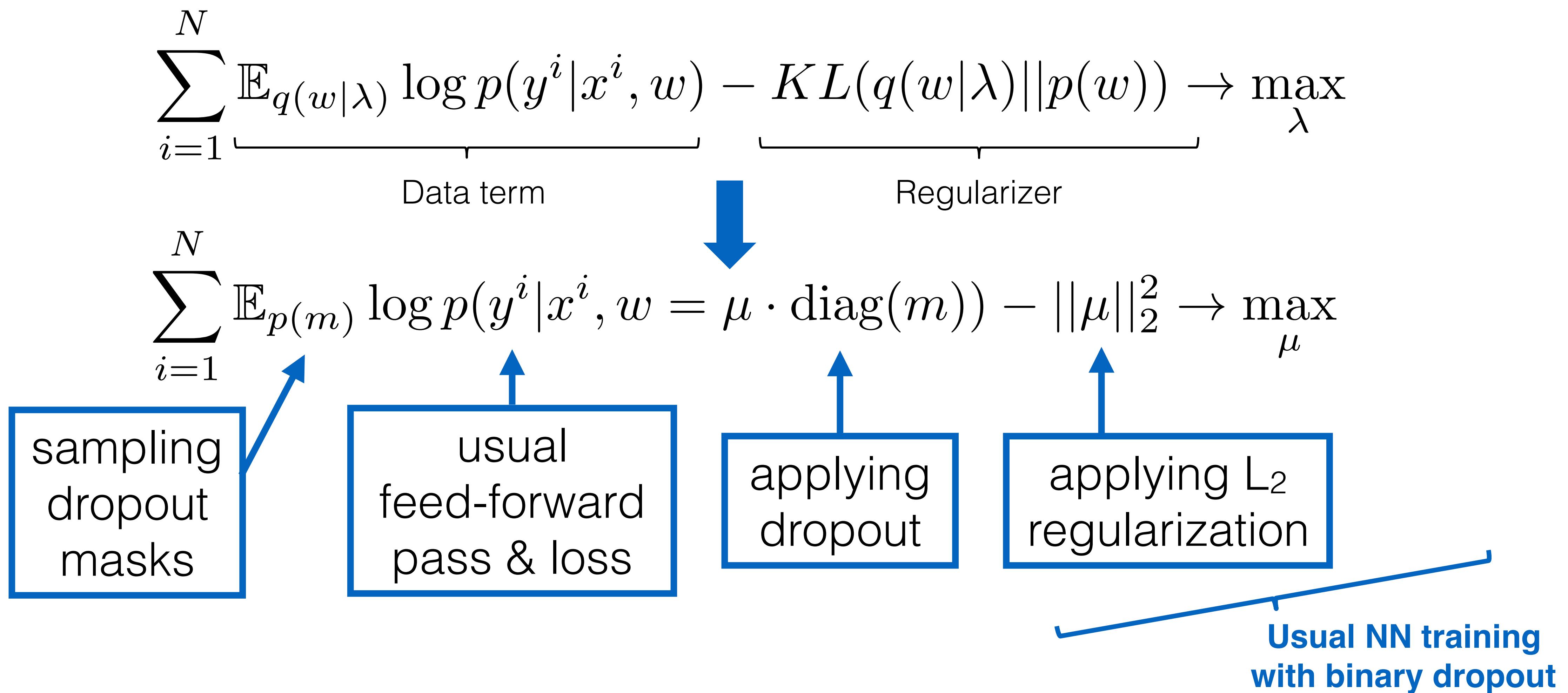
Approximate posterior: $q(W|\mu) = \prod_i q(w_i|\mu_i)$

$$q(w_i|\mu_i) = p \delta(0) + (1 - p) \delta(\mu_i)$$

Approximate KL-divergence: $KL(q(W|\mu)||p(W)) \approx \alpha \|\mu\|_2^2$

L2-regularization

Example: binary dropout, putting all together



Example: binary dropout, Bayesian benefits

Key messages of this example:

- Using binary dropout means being Bayesian!
- There are other dropout profits beyond regularization:
 - ensembling
 - uncertainty estimation

Plan

- Advantages of using Bayesian neural networks
- Training Bayesian neural networks
- First example: binary dropout
- Second example: Bayesian sparsification

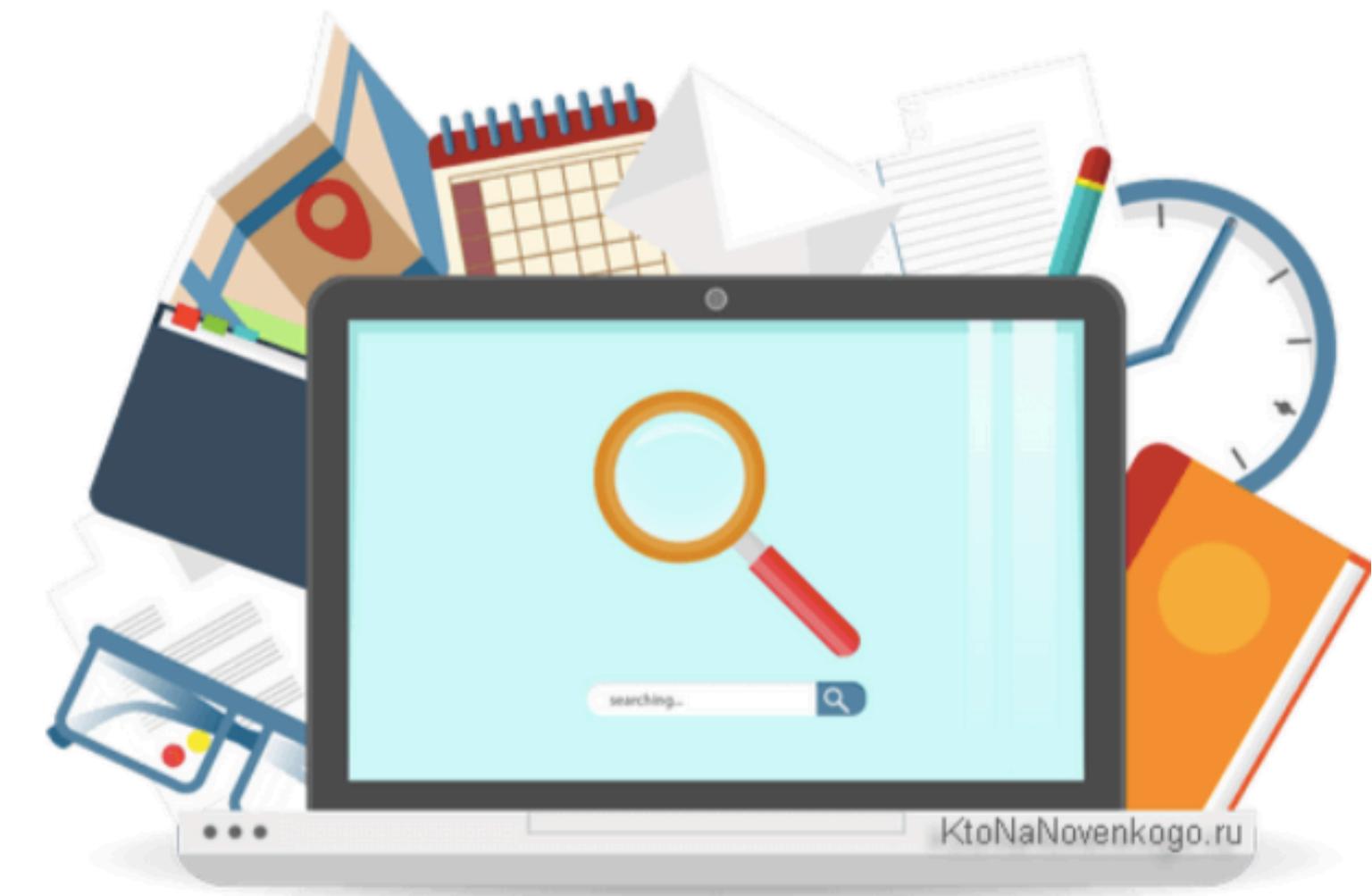
Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**



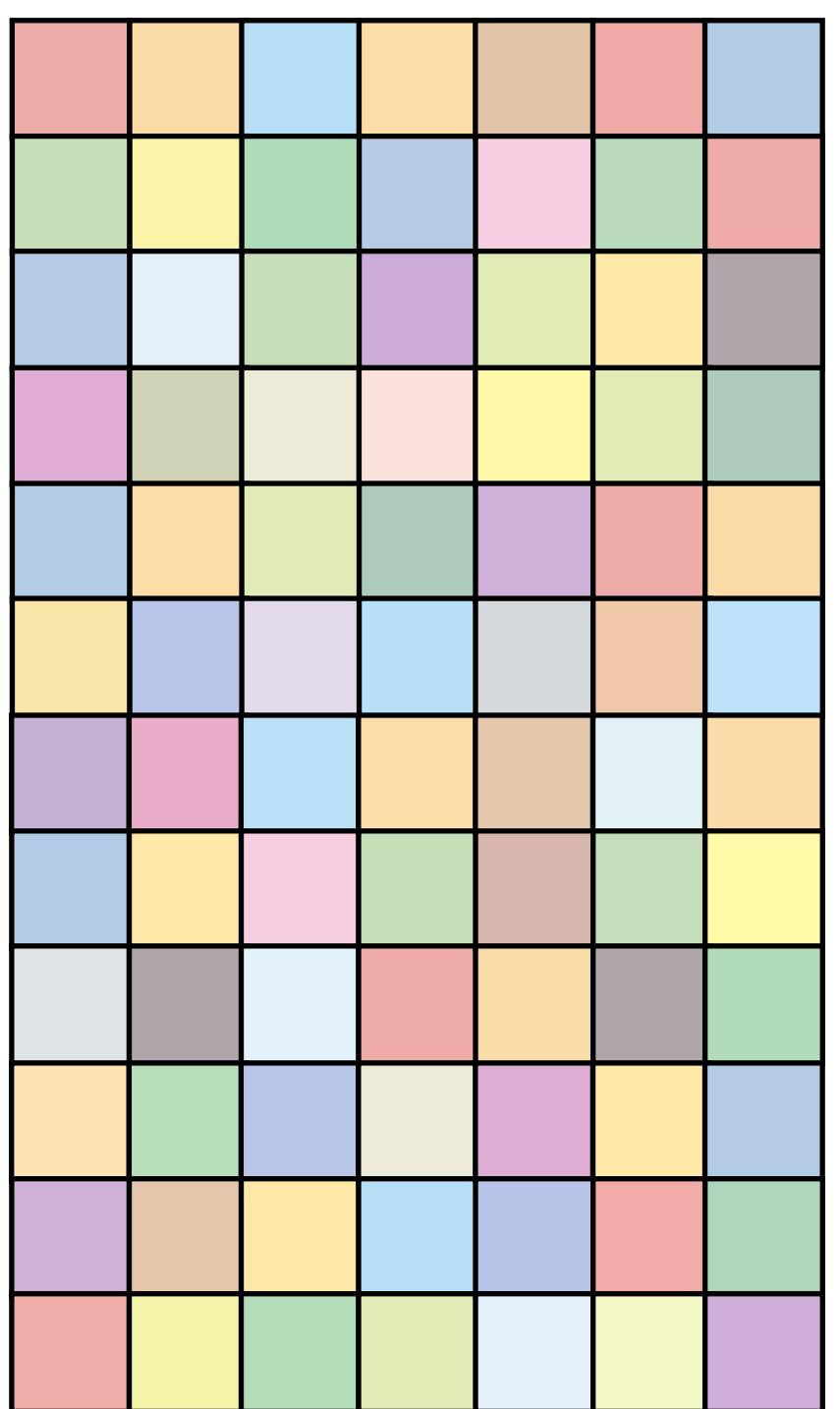
Compression of neural networks

- Deep neural networks achieve state-of-the-art performance in a variety of domains
- Model quality scales with model and dataset size
- State-of-the-art models usually incorporate **tens of millions of parameters**
- But **resources** (memory, processing time) **may be limited**

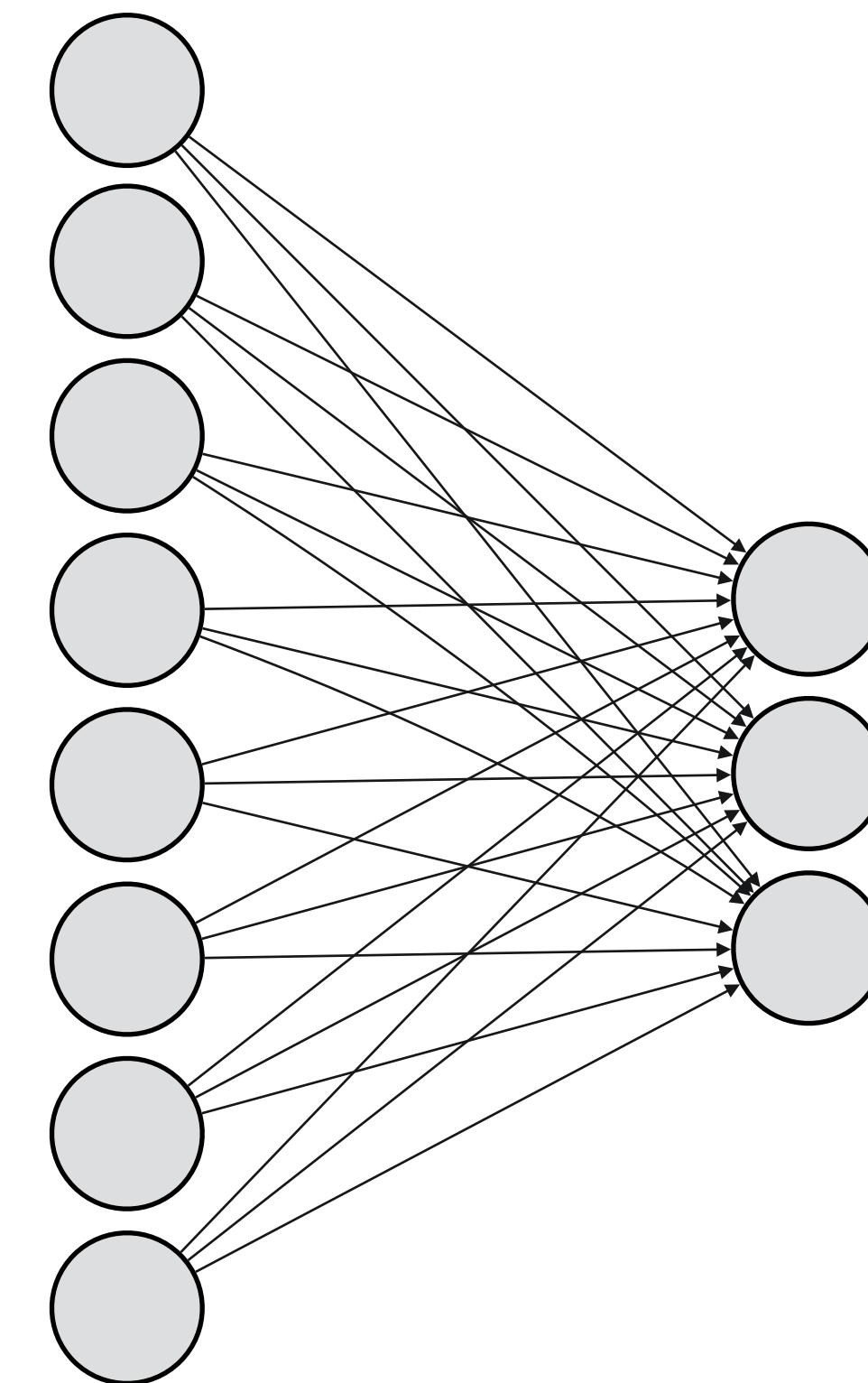


- One of the solutions — sparsification

Neural network

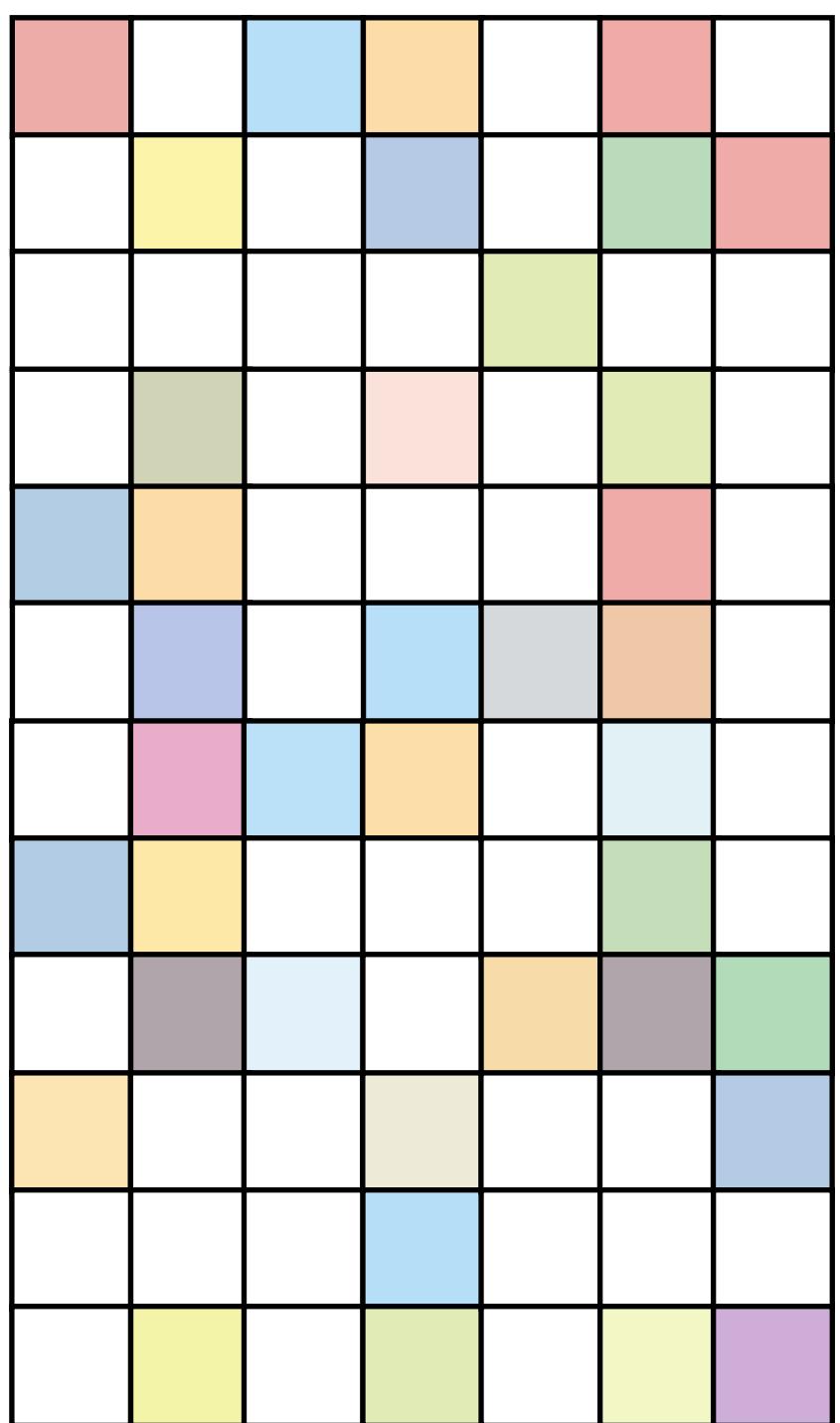


Weight matrix W



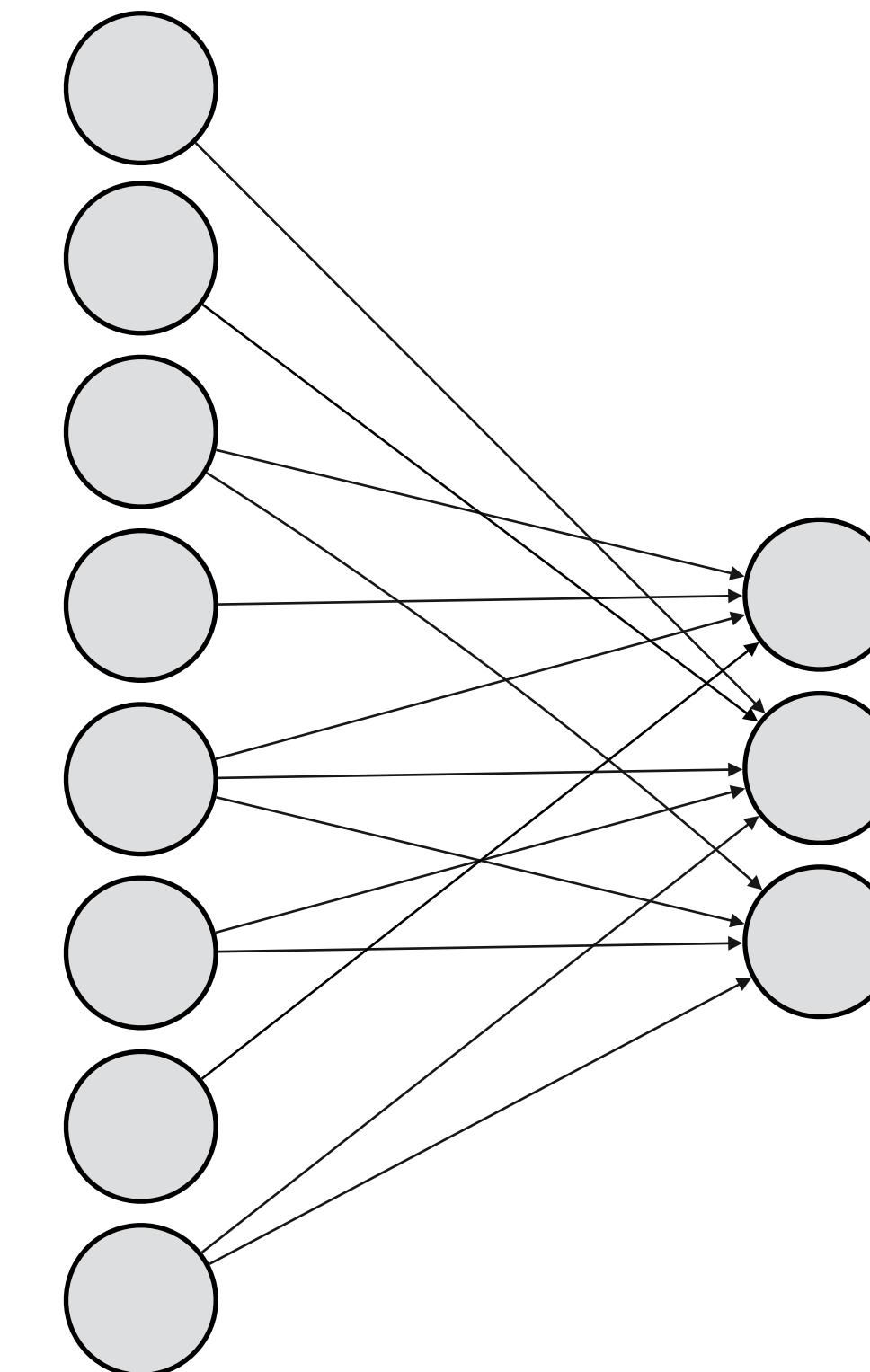
Computational graph

Sparse neural network



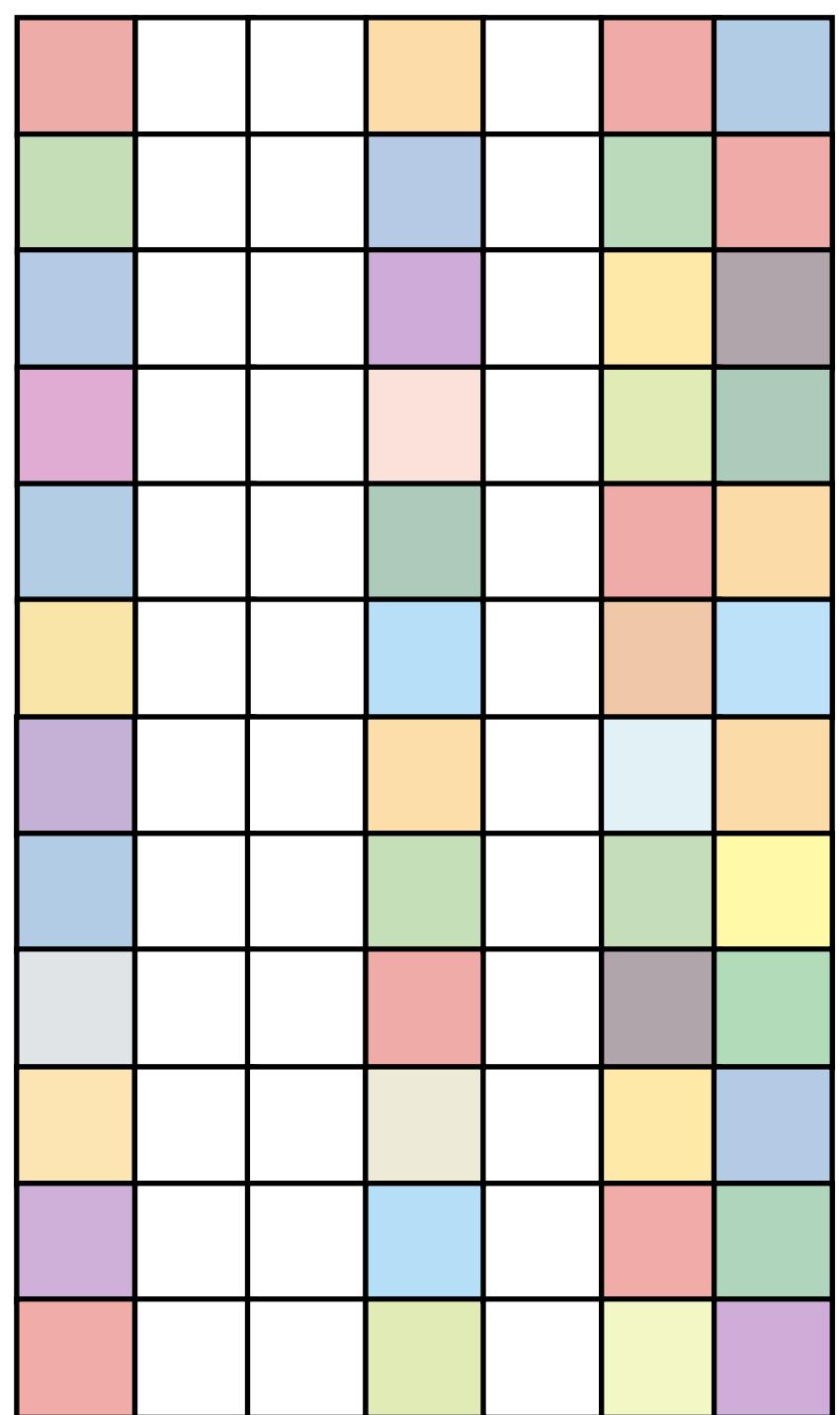
Weight matrix W

A lot of weights
set to zero



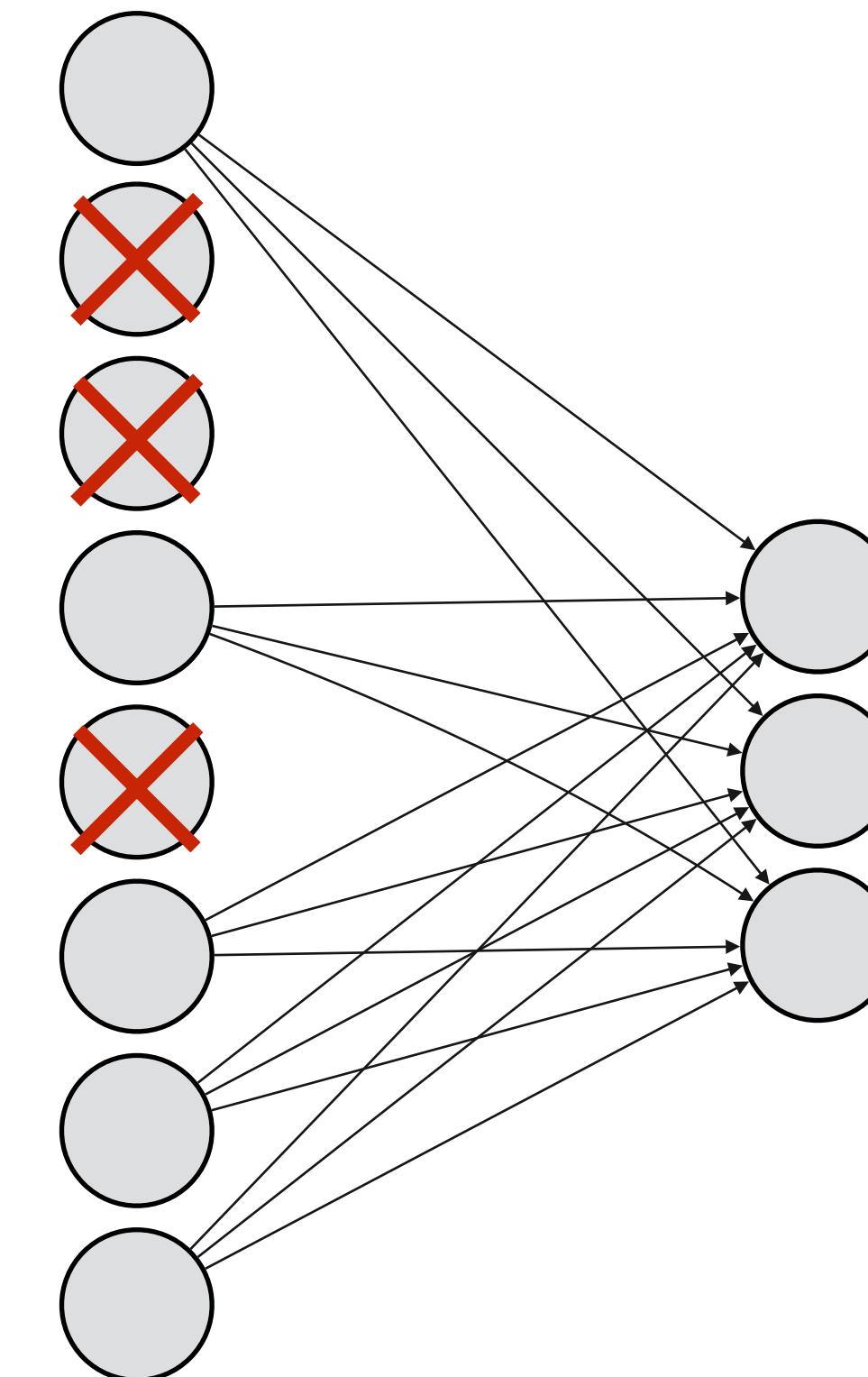
Computational graph

Structured sparsity



Weight matrix W

No outgoing edges
⇒ remove neuron



Computational graph

From general framework to particular method

$$\sum_{i=1}^N \mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w) - KL(q(w|\lambda)||p(w)) \rightarrow \max_{\lambda}$$

Model specification:

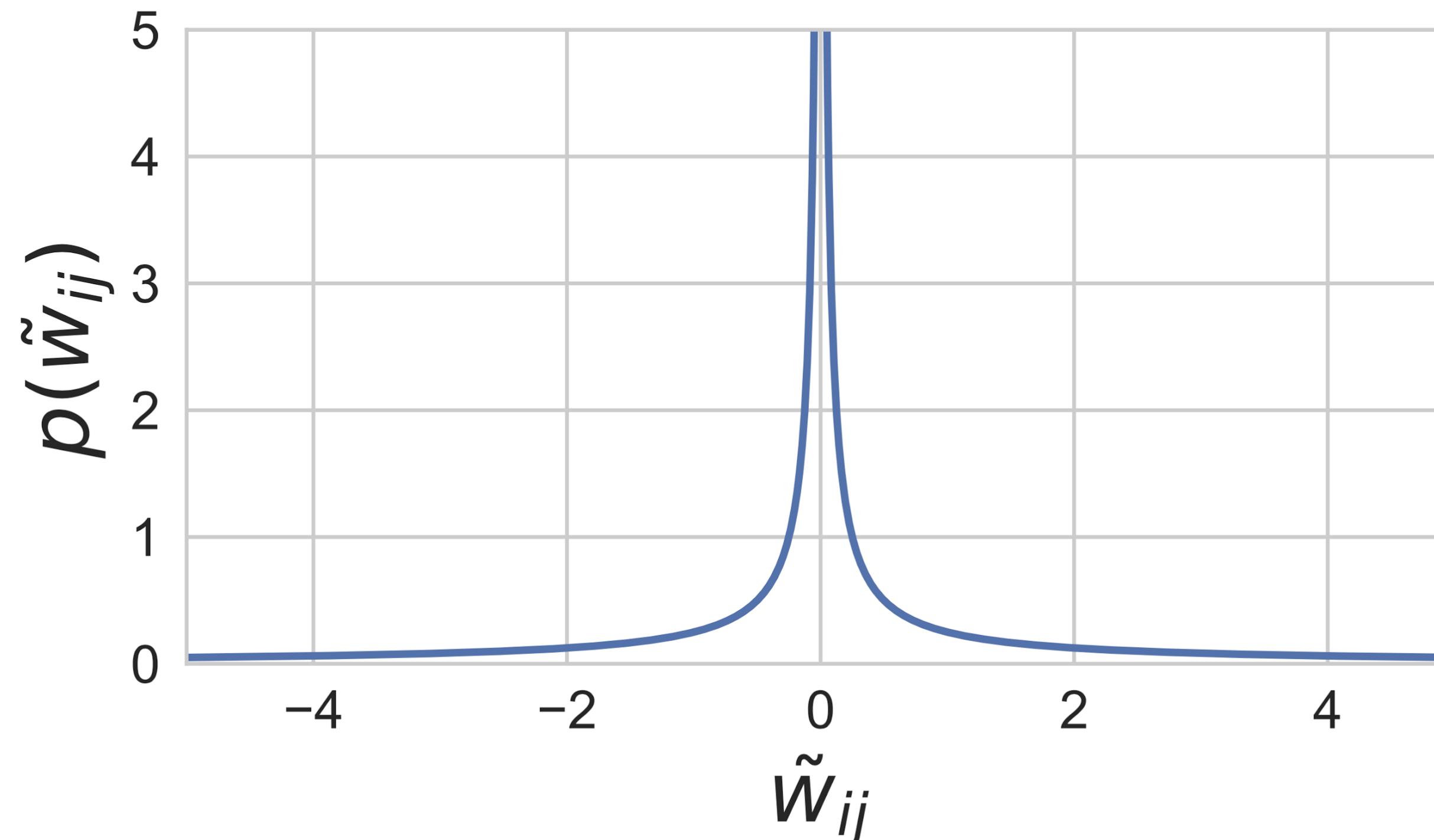
- Choose particular prior

Training:

- Choose particular family for approximate posterior
- How to compute the KL-divergence?

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$



Favors removing noisy weights!

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$

Approximate posterior: ?

Approximate KL-divergence: ?

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$

Approximate posterior: $q(w_{ij} | \mu_{ij}, \sigma_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$

Approximate KL-divergence: ?

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$

Approximate posterior: $q(w_{ij} | \mu_{ij}, \sigma_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$

Reparametrization: $\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$

Approximate KL-divergence: ?

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$

Approximate posterior: $q(w_{ij} | \mu_{ij}, \sigma_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$

Approximate KL-divergence: $-KL(q(w_{ij} | \mu_{ij}, \sigma_{ij}) \| p(w_{ij})) \approx f_{KL}(\alpha_{ij})$

$$\alpha_{ij} = \frac{\sigma_{ij}^2}{\mu_{ij}^2}$$

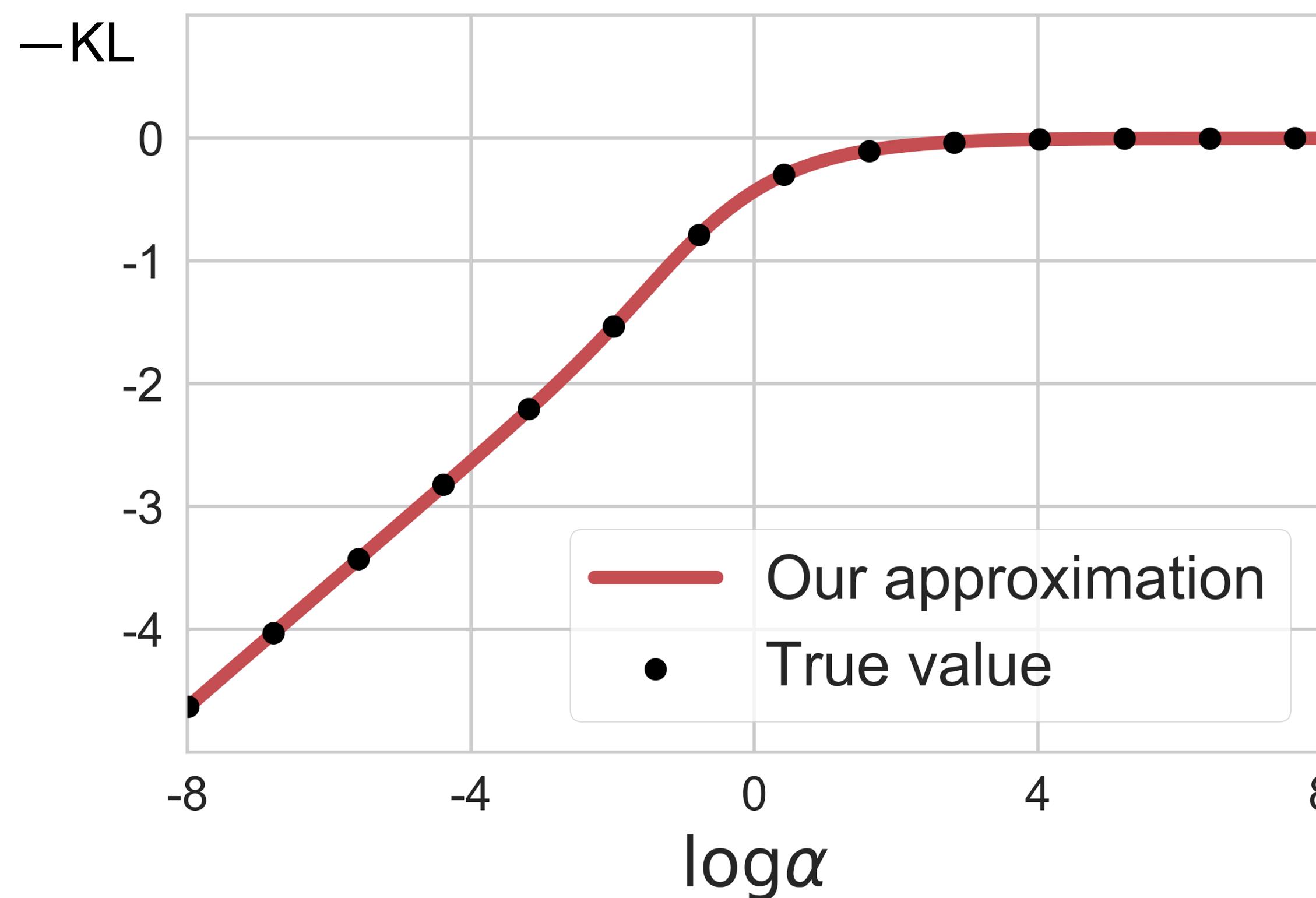
Favors large $\alpha_{ij} \Rightarrow$ removing noisy weights

Approximating KL-divergence

Remember: training Bayesian neural networks — optimizing ELBO:

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\mu, \sigma)} \log p(y^i|x^i, w)}_{\substack{\text{Data term} \\ \text{Sample weights from } q}} - \underbrace{KL(q(w|\mu, \sigma) || p(w))}_{\substack{\text{Regularizer} \\ \text{Analytical approximation}}} \rightarrow \max_{\mu, \log \sigma}$$

Approximating KL-divergence (fully factorized)



$$\begin{aligned}-\text{KL}(q(w_{ij}|\mu_{ij}, \sigma_{ij}) \| p(w_{ij})) &\approx \\ &\approx k_1 \sigma(k_2 + k_3 \log \alpha_{ij}) - 0.5 \log(1 + \alpha_{ij}^{-1}) + C\end{aligned}$$
$$k_1 = 0.63576 \quad k_2 = 1.87320 \quad k_3 = 1.48695$$

$$\alpha_{ij} = \frac{\sigma_{ij}^2}{\mu_{ij}^2}$$

- KL depends only on α_{ij}
- Favors large $\alpha_{ij} \Rightarrow$ removing noisy weights

Ok, sparsify weights. What about biases?

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma}$$

Treat biases as deterministic parameters and find a point estimate:

$$\sum_{i=1}^N \mathbb{E}_{q(w|\mu,\sigma)} \log p(y^i|x^i, w, \mathbf{b}) - KL(q(w|\mu,\sigma)||p(w)) \rightarrow \max_{\mu, \log \sigma, \mathbf{b}}$$

Example: sparse variational dropout

Prior: $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$

Approximate posterior: $q(w_{ij} | \mu_{ij}, \sigma_{ij}) = \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$

Approximate KL-divergence: $-KL(q(w_{ij} | \mu_{ij}, \sigma_{ij}) \| p(w_{ij})) \approx f_{KL}(\alpha_{ij})$

$$\alpha_{ij} = \frac{\sigma_{ij}^2}{\mu_{ij}^2}$$

Favors large $\alpha_{ij} \Rightarrow$ removing noisy weights

Final algorithm

Training on a mini-batch X with labels Y :

1. Sample weights: $\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}, \quad \hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$
2. Forward pass: $Y_{\text{pred}} = NN(X, \hat{w}, b)$
3. Backward pass + SGD step: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{LOSS}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Final algorithm

Training on a mini-batch X with labels Y :

1. Sample weights: $\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}$, $\hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$
2. Forward pass: $Y_{\text{pred}} = NN(X, \hat{w}, b)$
3. Backward pass + SGD step: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Pruning after training:

If $\mu_{ij}^2/\sigma_{ij}^2 < \text{threshold}$:

$$\mu_{ij} = 0, \sigma_{ij} = 0$$

signal-to-noise ratio

Final algorithm

Training on a mini-batch X with labels Y :

1. Sample weights: $\hat{w}_{ij} = \mu_{ij} + \hat{\epsilon}_{ij}\sigma_{ij}$, $\hat{\epsilon}_{ij} \sim \mathcal{N}(0, 1)$
2. Forward pass: $Y_{\text{pred}} = NN(X, \hat{w}, b)$
3. Backward pass + SGD step: compute stochastic gradients of ELBO:

$$\nabla_{\mu, \log \sigma, b} \left(N \cdot \text{Loss}(Y, Y_{\text{pred}}) + \text{SparseReg}(\sigma/\mu) \right)$$

Pruning after training:

If $\mu_{ij}^2/\sigma_{ij}^2 <$ threshold:

$$\mu_{ij} = 0, \sigma_{ij} = 0$$

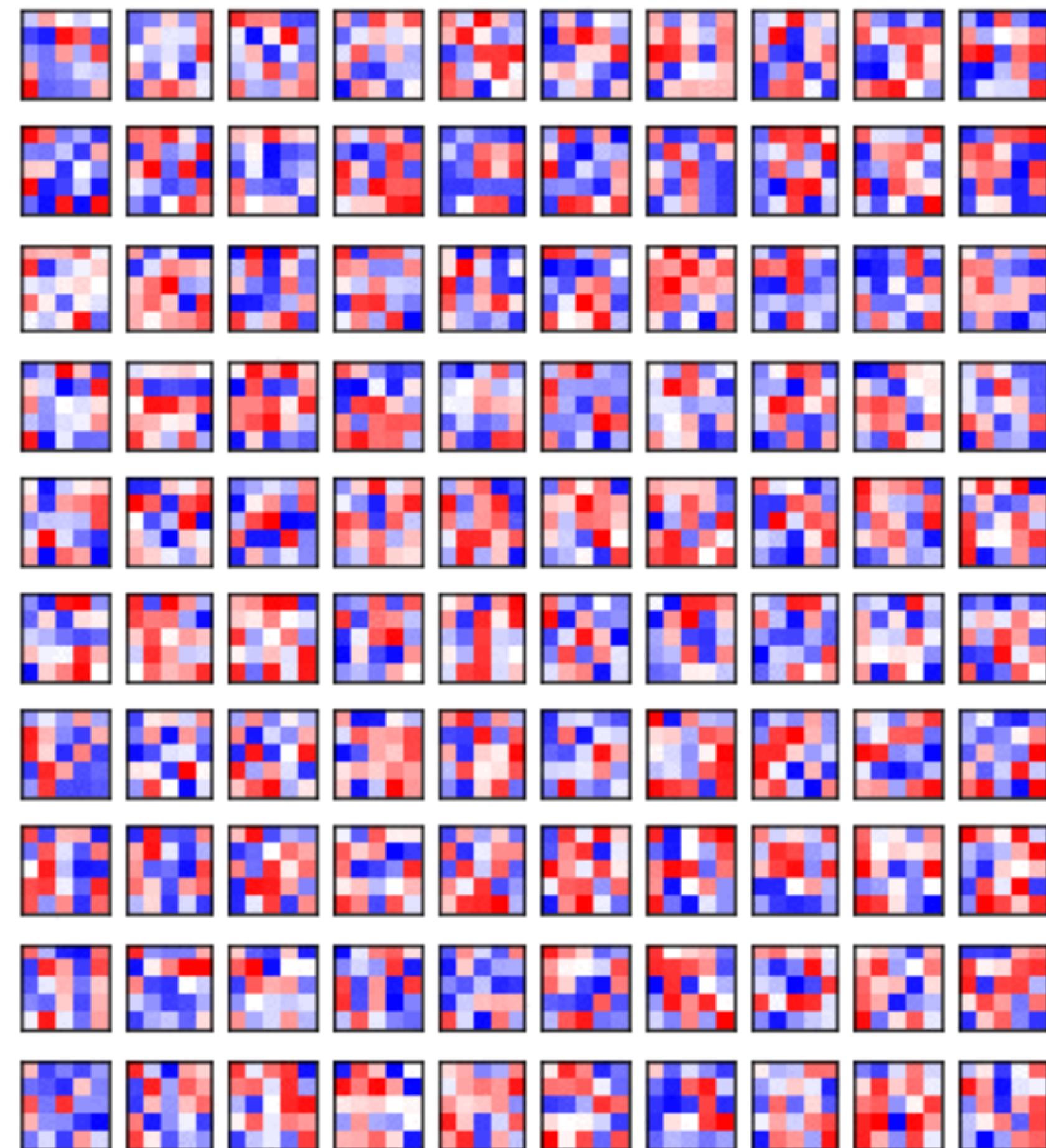
Prediction for a mini-batch X :

Return $Y_{\text{pred}} = NN(X, \mu, b)$

do not ensemble because we want
the most compact and fast network

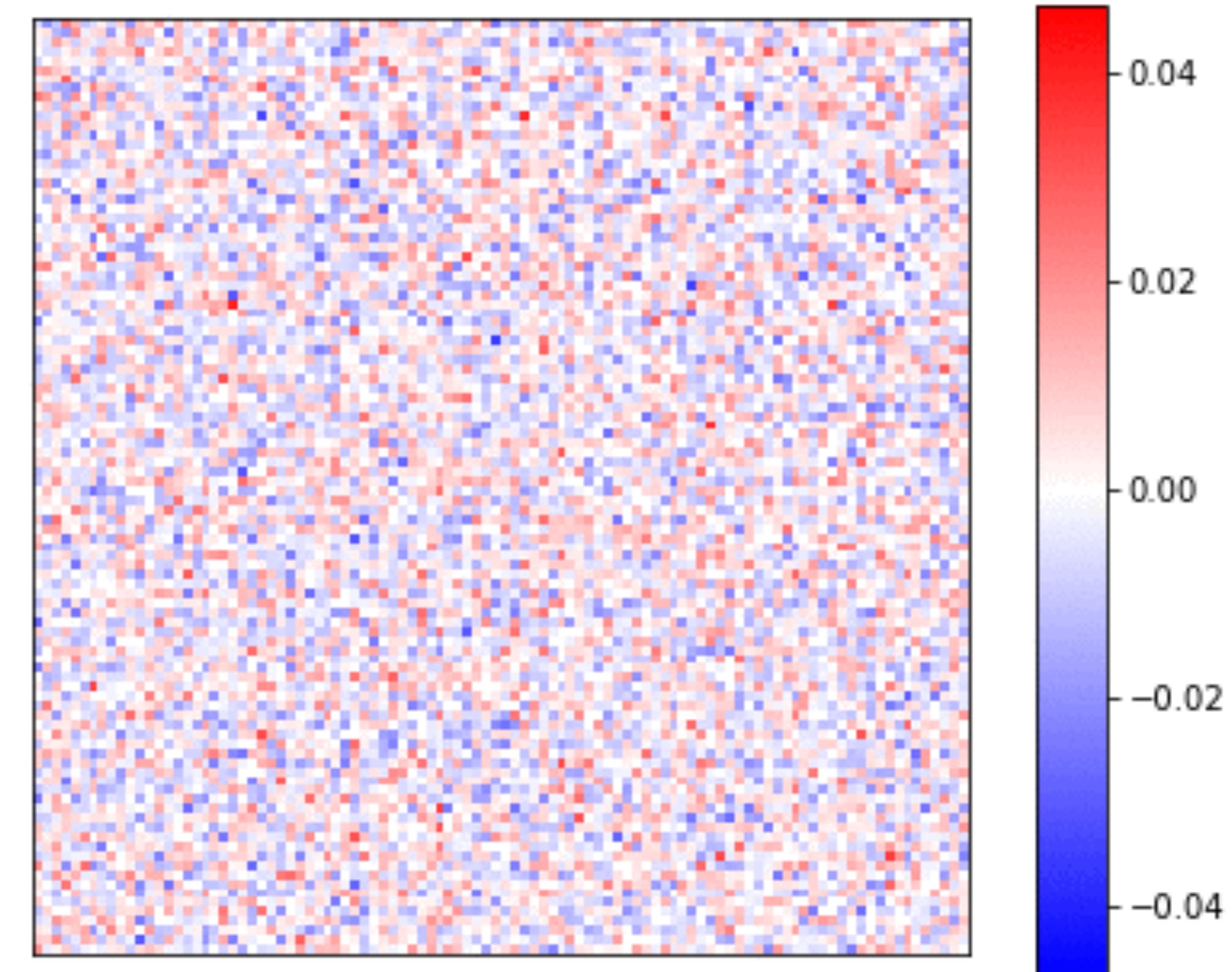
Sparse variational dropout: visualization

Epoch: 0 Compression ratio: 1x Accuracy: 8.4



LeNet-5: convolutional layer

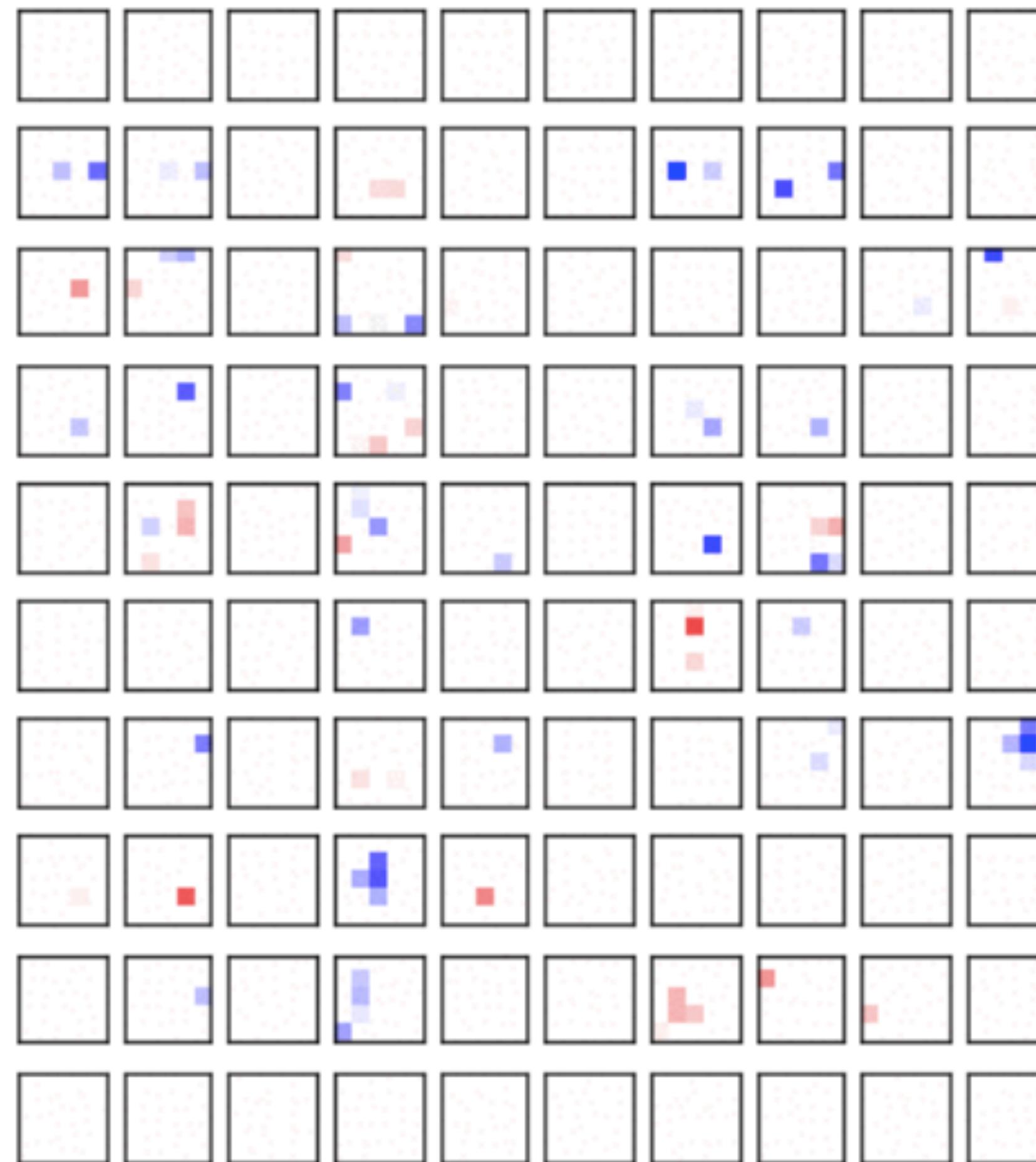
Epoch: 0 Compression ratio: 1x Accuracy: 8.4



LeNet-5: fully-connected layer
(100 x 100 patch)

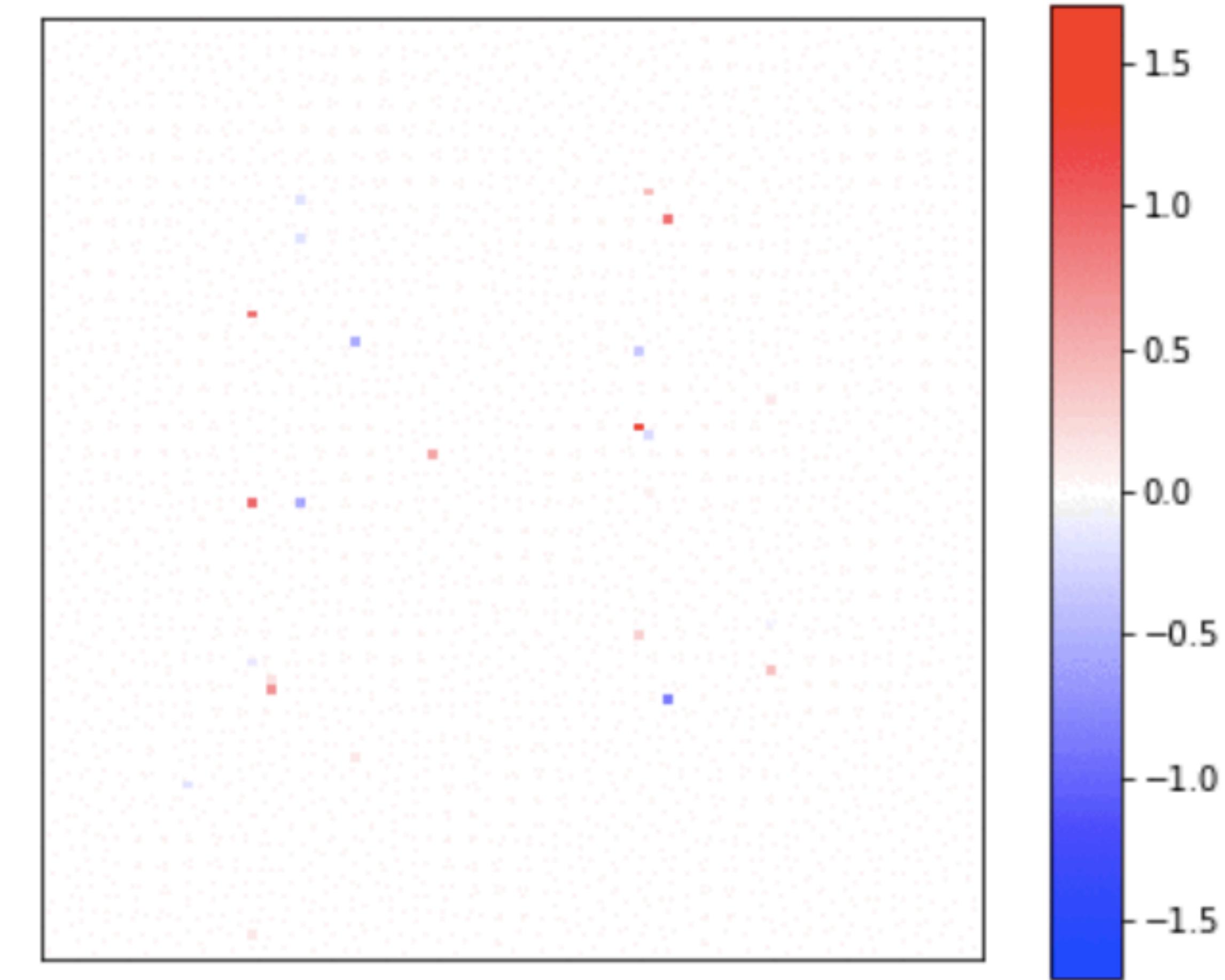
Sparse variational dropout: visualization

Epoch: 200 Compression ratio: 270x Accuracy: 99.3



LeNet-5: convolutional layer

Epoch: 200 Compression ratio: 270x Accuracy: 99.3

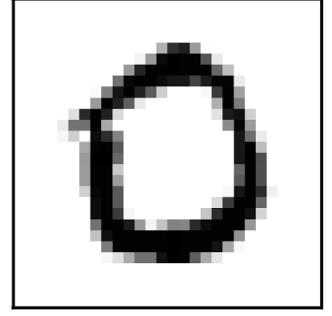
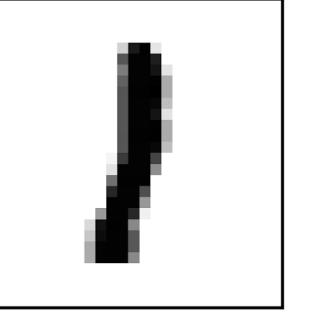
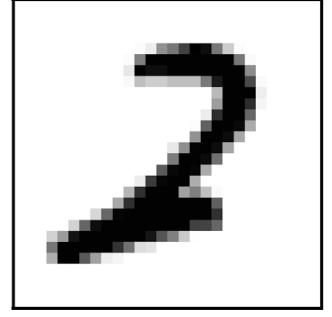
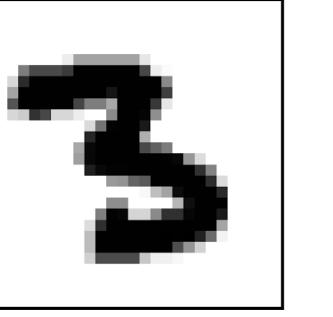
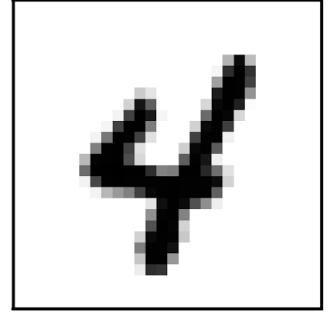
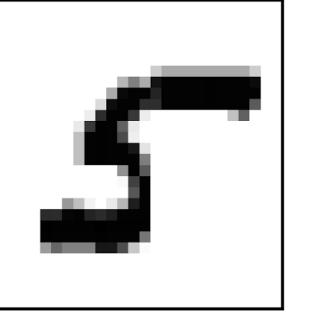
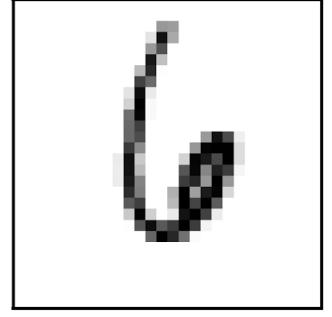
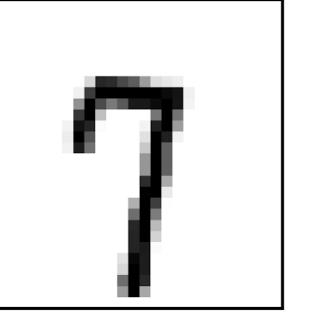
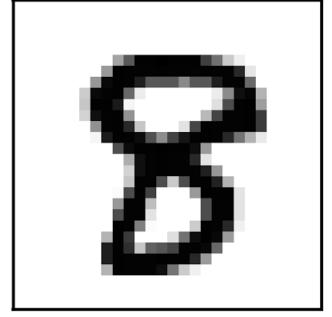
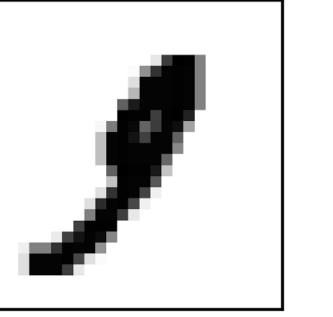


LeNet-5: fully-connected layer
(100 x 100 patch)

Lenet-5-Caffe and Lenet-300-100 on MNIST

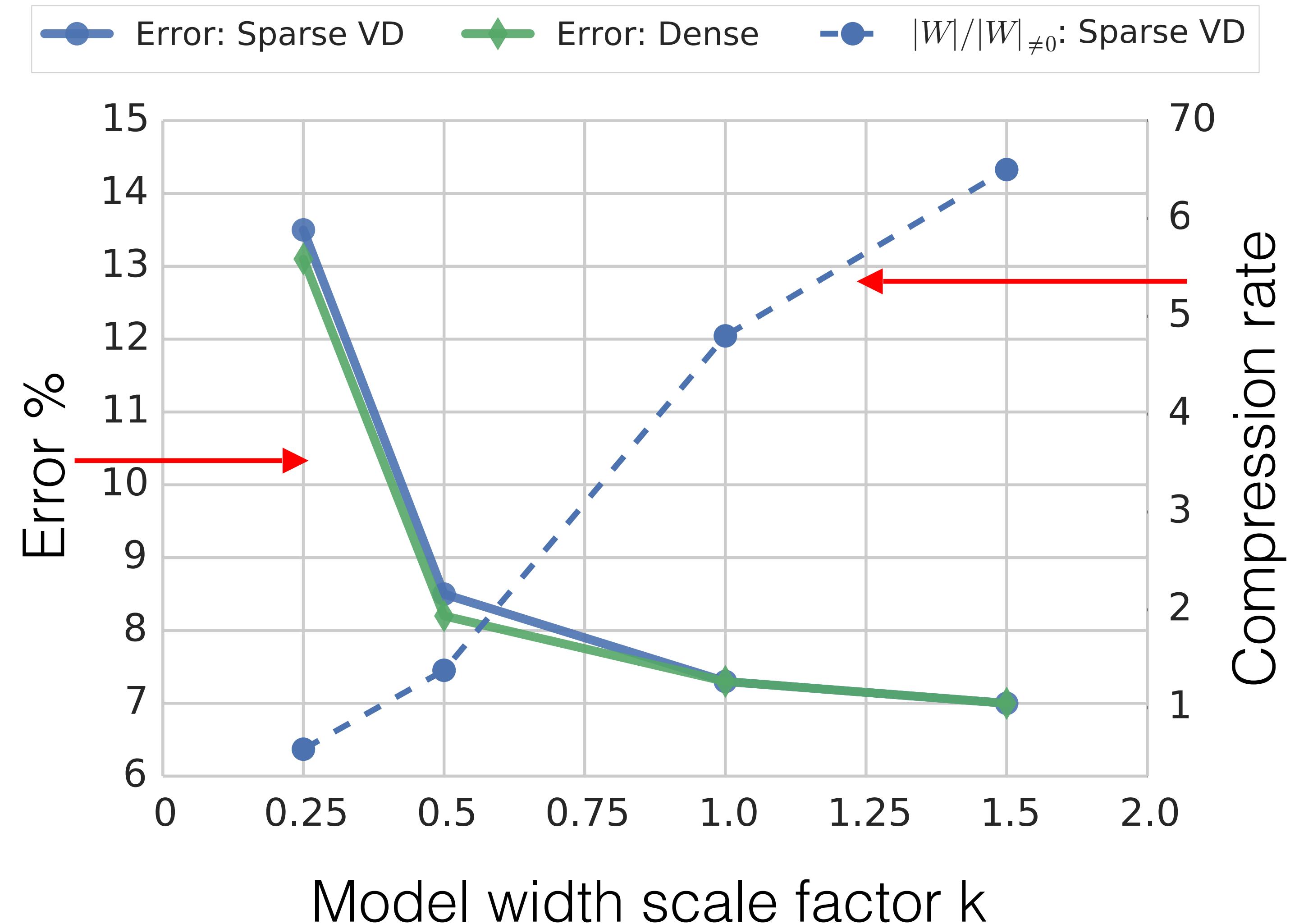
Fully Connected network: LeNet-300-100

Convolutional network: Lenet-5-Caffe

Network	Method	Error %	Sparsity per Layer %	$\frac{ \mathbf{W} }{ \mathbf{W}_{\neq 0} }$		
LeNet-300-100	Original	1.64		1		
	Pruning	1.59	92.0 – 91.0 – 74.0	12		
	DNS	1.99	98.2 – 98.2 – 94.5	56		
	SWS	1.94		23		
	(ours) Sparse VD	1.92	98.9 – 97.2 – 62.0	68		
LeNet-5-Caffe	Original	0.80		1		
	Pruning	0.77	34 – 88 – 92.0 – 81	12		
	DNS	0.91	86 – 97 – 99.3 – 96	111		
	SWS	0.97		200		
	(ours) Sparse VD	0.75	67 – 98 – 99.8 – 95	280		

VGG-like on CIFAR-10

Number of filters / neurons is linearly scaled by k (the width of the network)



Random Labeling



Dataset	Architecture	Train Acc.	Test Acc.	Sparsity
MNIST	FC + BD	100%	10%	—
MNIST	FC + Sparse VD	10%	10%	100%
CIFAR-10	VGG + BD	100%	10%	—
CIFAR-10	VGG + Sparse VD	10%	10%	100%

No dependency between data and labels \Rightarrow Sparse VD yields an empty model
where conventional models easily overfit.

Sparse variational dropout: key messages

- Prior distribution can encode our desirable model properties (e. g. sparse weights)
- Other Bayesian compression techniques:
 - group sparsification (removing neurons / filters)
 - quantization (low-precision weights)

Summary

- A lot of BNN advantages: regularization, ensembling, uncertainty estimation, ...
- To train BNN, one should optimize ELBO using DSVI & RT
- Three steps towards a particular method
- Using binary dropout means being Bayesian
- Prior distribution can encode our desirable model properties