

Bayesian neural networks

(and VI in implicit models)

Dmitry Molchanov

Samsung AI Center, Samsung-HSE Laboratory



Lecture outline

- What are Bayesian Neural Networks (BNNs)
- Why go Bayesian
- How to train BNNs
- Variational inference with implicit posteriors

What you already know

- Stochastic optimization
- Bayesian modelling
- Latent variable models
- Variational inference
 - Bayesian inference \leftrightarrow (stochastic) optimization
 - (Doubly) Stochastic variational inference
 - Reparameterization Trick



\Rightarrow **Bayesian neural networks**

Regularization by noise

Traditional (1943+) regularization: add some penalty for model complexity

- Norm-based regularization (L2, L1)

$$Objective = DataLoss(X, W) + Regularizer(W)$$

- Max norm constraint

More recent (1990+) approaches: regularization by noise

- Input noise:
 - Denoising autoencoders
 - Data augmentation
- Weight noise:
 - Dropout (reviving noise regularization in 2012)
 - Gaussian weight noise

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

- Gradient noise

Bayesian framework provides a principled approach to training with noise!

Generative models vs discriminative models

Bayesian Discriminative Model:

Likelihood $p(\mathbf{t}|X, \mathbf{w}) = \prod_{i=1}^N p(t_i|\mathbf{x}_i, \mathbf{w})$ Can be a neural network with weights W !

Prior $p(\mathbf{w})$

Posterior $p(\mathbf{w}|X, \mathbf{t}) = \frac{p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})d\mathbf{w}} = ?$

- No local latent variables; we want the posterior over the weights instead
- Much higher dimensionality
 - 10^2 - 10^3 for generative models, 10^5 - 10^8 and more for discriminative models

Why go Bayesian?

A principled framework with many useful applications

- Regularization
- Ensembling
- Uncertainty estimation
- On-line / continual learning
- And more (stay tuned for the next lecture!)

Ensembling

A Bayesian neural network is **an infinite ensemble** of neural networks

$\mathbf{w} \sim p(\mathbf{w}|X, \mathbf{t})$ One sample from the posterior
One element of the ensemble

Predictive distribution $p(t^*|\mathbf{x}^*, X, \mathbf{t}) = \int p(t^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|X, \mathbf{t})d\mathbf{w}$

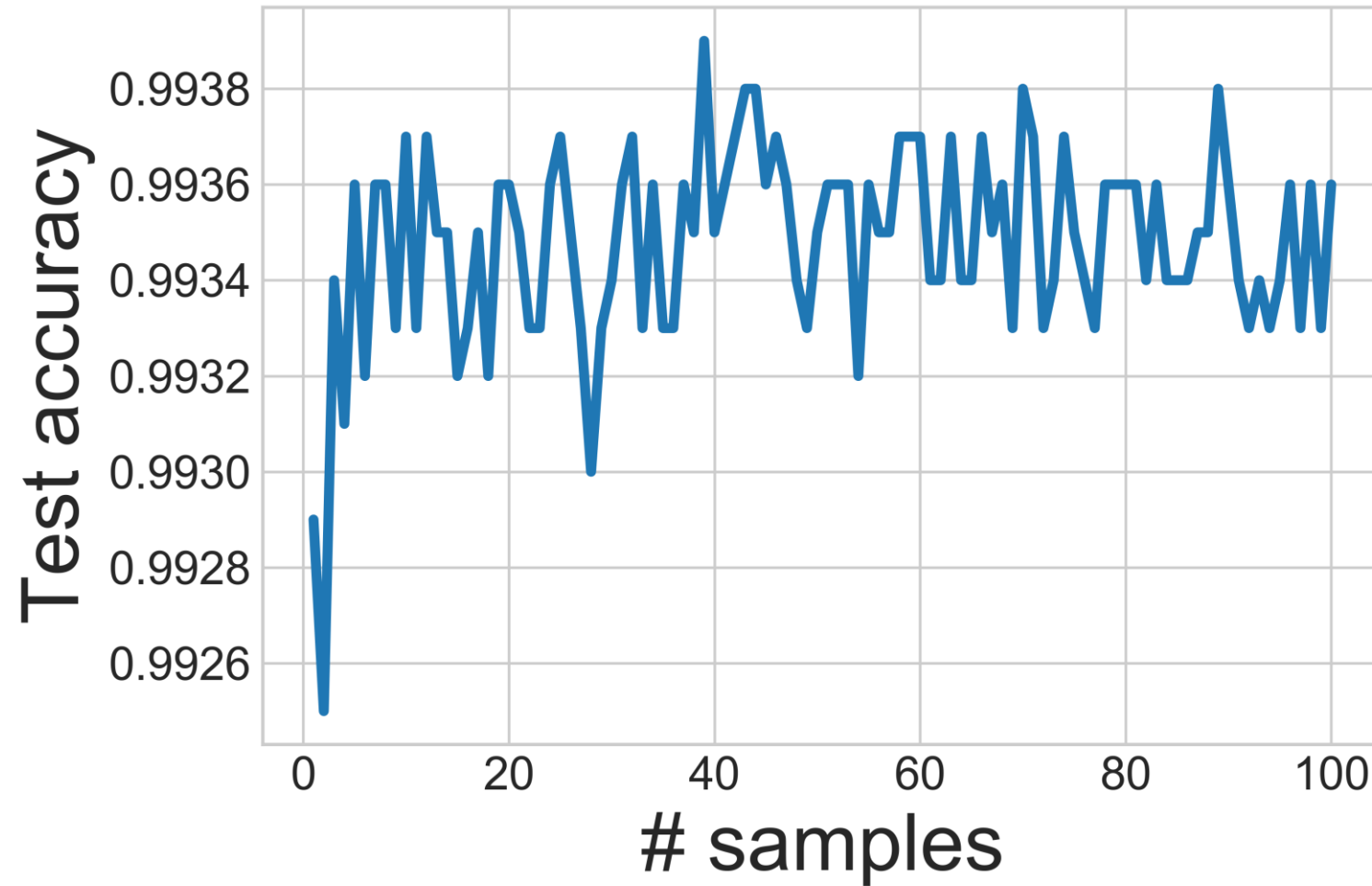
And its unbiased estimate

$$\mathbb{E}_{p(\mathbf{w}|X, \mathbf{t})}p(t^*|\mathbf{x}^*, \mathbf{w}) \simeq \frac{1}{K} \sum_{i=1}^K p(t^*|\mathbf{x}^*, \mathbf{w}^k); \quad \mathbf{w}^k \sim p(\mathbf{w}|X, \mathbf{t})$$

- Higher accuracy
- More robust

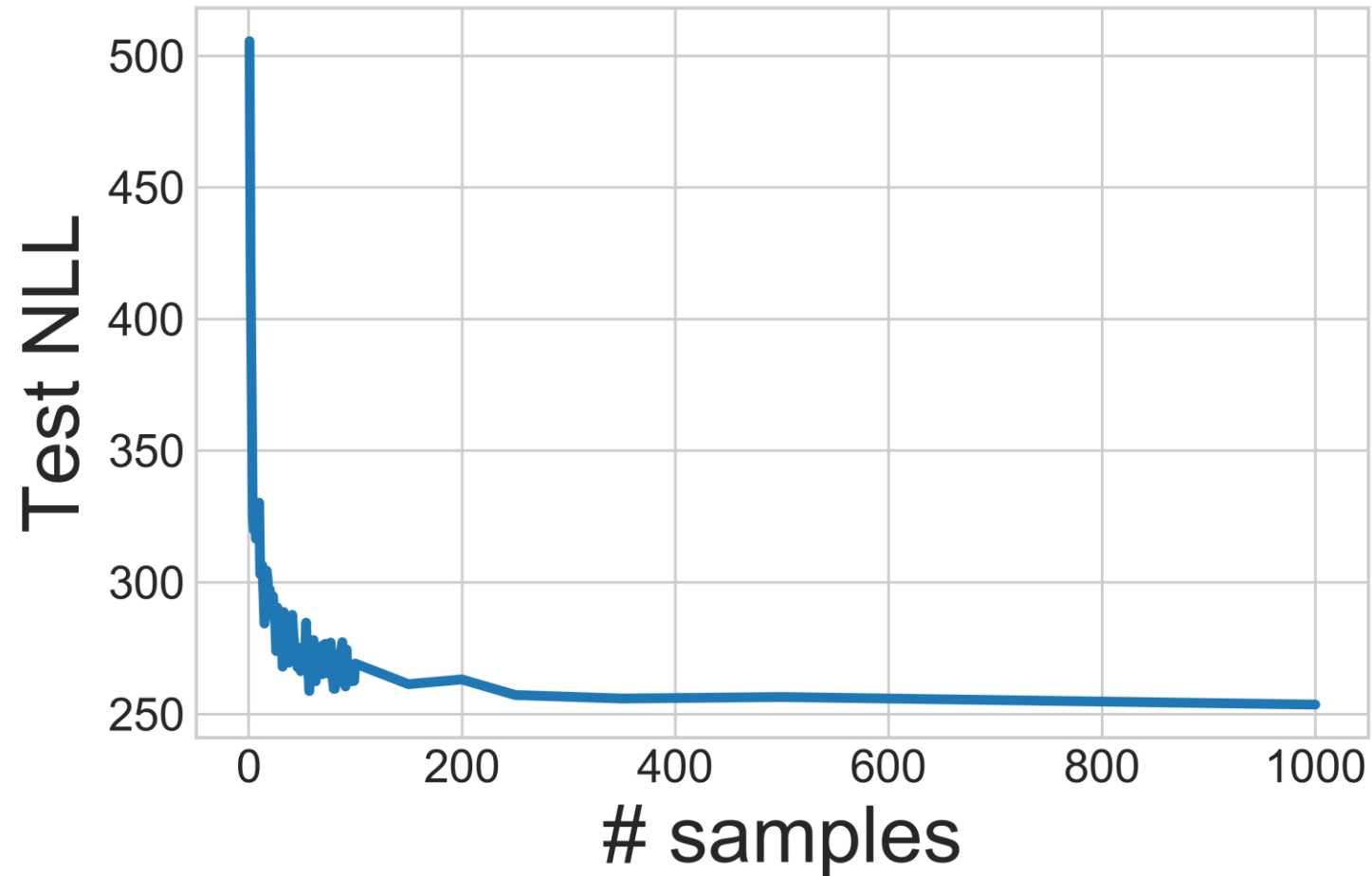
Average SoftMax outputs
across several samples

Ensembling



Accuracy quickly
saturates

Ensembling

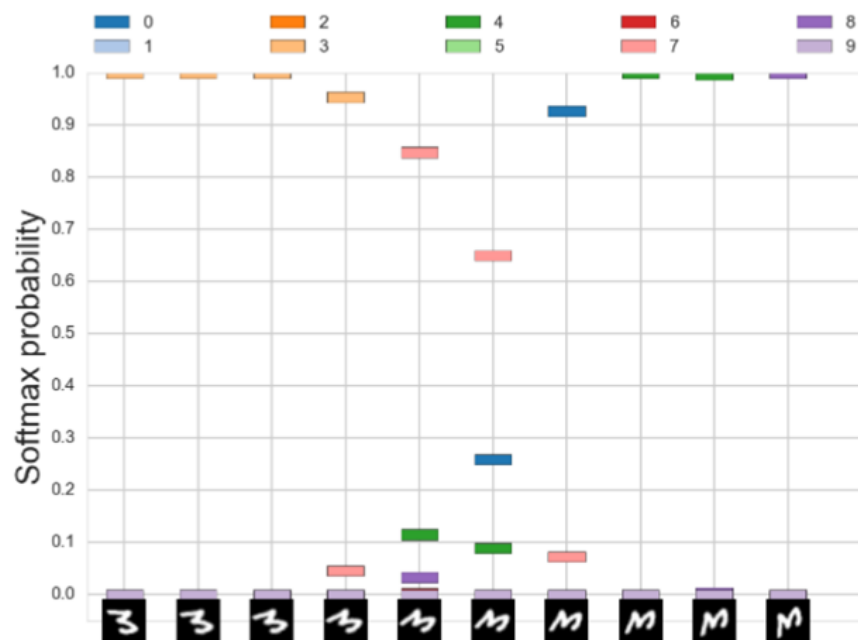


But the NLL keeps improving!
This is a measure of “uncertainty”

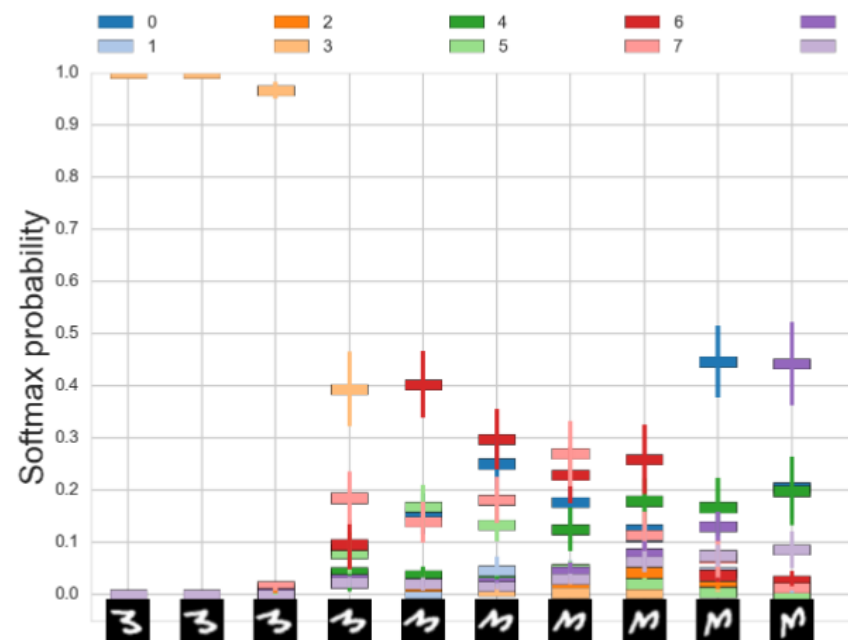
Uncertainty estimation

Deterministic NNs: a **point estimate** of the output, overconfident

Bayesian framework allows us to obtain a **distribution** over the outputs



(a) LeNet with weight decay



(b) LeNet with multiplicative formalizing flows

Model selection and compression

- Empirical Bayes (maximum evidence)
 - Choose hyperparameters
 - Model compression
 - Similar to the Relevance Vector Machine
- Special sparsity-inducing priors
 - Stay tuned for the next lecture

On-line / incremental learning

Assume that the dataset arrives in independent parts.

$$\mathcal{D} = \mathcal{D}_1 \cup \mathcal{D}_2 \cup \dots \cup \mathcal{D}_M$$

We can train on the first dataset as usual...

$$p(\mathbf{w}|\mathcal{D}_1) = \frac{p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})}{\int p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})d\mathbf{w}}$$

... And then use the obtained posterior as the prior for the next step!

$$\begin{aligned} p(\mathbf{w}|\mathcal{D}_2, \mathcal{D}_1) &= \frac{p(\mathcal{D}_2|\mathbf{w})p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})}{\int p(\mathcal{D}_2|\mathbf{w})p(\mathcal{D}_1|\mathbf{w})p(\mathbf{w})d\mathbf{w}} = \\ &= \frac{p(\mathcal{D}_2|\mathbf{w})p(\mathbf{w}|\mathcal{D}_1)}{\int p(\mathcal{D}_2|\mathbf{w})p(\mathbf{w}|\mathcal{D}_1)d\mathbf{w}} \end{aligned}$$

Using these sequential updates, we can find $p(\mathbf{w}|\mathcal{D})$!

Variational inference for Bayesian NNs

The posterior distribution $p(\mathbf{w}|X, \mathbf{t}) = \frac{p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{t}|X, \mathbf{w})p(\mathbf{w})d\mathbf{w}}$

How to find it? Use (doubly stochastic) variational inference!

$$\begin{aligned} q(\mathbf{w}|\boldsymbol{\phi}) &\approx p(\mathbf{w}|X, \mathbf{t}) \\ \text{KL}(q(\mathbf{w}|\boldsymbol{\phi}) \parallel p(\mathbf{w}|X, \mathbf{t})) &\rightarrow \min_{\boldsymbol{\phi}} \\ \mathcal{L}(\boldsymbol{\phi}) = \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\phi})} \log p(\mathbf{t}|X, \mathbf{w}) - \text{KL}(q(\mathbf{w}|\boldsymbol{\phi}) \parallel p(\mathbf{w})) &\rightarrow \max_{\boldsymbol{\phi}} \end{aligned}$$

Only two differences from LVMs:

- KL-term is global
- Extremely high-dimensional posterior

Reparameterization trick for Bayesian NNs

Reparameterize $q(\mathbf{w}|\boldsymbol{\phi})$ and plug the sample into the ELBO

$$\begin{aligned} \mathbf{w} \sim q(\mathbf{w}|\boldsymbol{\phi}) &\Leftrightarrow \mathbf{w} = g(\boldsymbol{\epsilon}, \boldsymbol{\phi}); \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon}) \\ \mathcal{L}(\boldsymbol{\phi}) &= \mathbb{E}_{p(\boldsymbol{\epsilon})} \log p(\mathbf{t}|X, \mathbf{w} = g(\boldsymbol{\epsilon}, \boldsymbol{\phi})) - \text{KL}(q \parallel p) \rightarrow \max_{\boldsymbol{\phi}} \end{aligned}$$

Obtain an unbiased differentiable mini-batch estimator

$$\mathcal{L}(\boldsymbol{\phi}) \simeq \sum_i \log p(\mathbf{t}_{m_i} | \mathbf{x}_{m_i}, \mathbf{w} = g(\boldsymbol{\epsilon}, \boldsymbol{\phi})) - \text{KL}(q \parallel p); \quad \boldsymbol{\epsilon} \sim p(\boldsymbol{\epsilon})$$

Very similar to conventional loss functions

Basically, using any kind of noise during training is close to being Bayesian

Usually just **1 sample per iteration is enough!**

Ex: dropout training as variational inference

Binary dropout results in a binary dropout posterior

$$\mathbf{w} = \boldsymbol{\phi} \cdot \text{diag}(\boldsymbol{\epsilon}); \quad \epsilon_i \sim \text{Bernoulli}(p)$$

It can be shown that a Gaussian prior leads to L2 regularization here
ELBO for binary dropout training:

$$\mathcal{L}(\boldsymbol{\phi}) = \mathbb{E}_{p(\boldsymbol{\epsilon})} \log p(\mathbf{t} | X, \mathbf{w} = \boldsymbol{\phi} \cdot \text{diag}(\boldsymbol{\epsilon})) - \lambda \|\boldsymbol{\phi}\|_2^2 \rightarrow \max_{\boldsymbol{\phi}}$$

- Using binary dropout means being Bayesian 😊
- There are other uses beyond regularization!
 - Ensembling
 - Uncertainty estimation
 - We can tune the dropout rate p using REINFORCE and extensions

Gal, Yarin, and Zoubin Ghahramani. "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning." *ICML* 2016.

Ex: Fully-Factorized Gaussians

Approximate posterior

$$q(\mathbf{w}) = \prod_i \mathcal{N}(w_i | \mu_i, \sigma_i^2)$$

Reparameterization

$$\mathbf{w} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}; \quad \epsilon_i \sim \mathcal{N}(0, 1)$$

The prior here is, e.g. a zero-centered FF Gaussian prior with variance σ_{prior}^2

ELBO:

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\sigma}) = \mathbb{E}_{p(\boldsymbol{\epsilon})} \log p(\mathbf{t} | X, \mathbf{w} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}) - \underbrace{\frac{\|\boldsymbol{\mu}\|_2^2 + \|\boldsymbol{\sigma}\|_2^2}{2\sigma_{prior}^2} + \sum_i \log \frac{\sigma_i^2}{\sigma_{prior}^2}}_{\text{KL between two } \mathcal{N}} \rightarrow \max_{\boldsymbol{\mu}, \boldsymbol{\sigma}}$$


- More tractable
- Richer approximation
- Twice as many parameters
- Start with small σ , optimize w.r.t. $\log \sigma$ to avoid constrained optimization

The local reparameterization trick

ELBO estimator may have high variance:

$$\mathcal{L}(\phi) \simeq \hat{\mathcal{L}}(\phi) = \frac{N}{M} \sum_{i=1}^M L_i(\phi, \epsilon)$$

Shared noise sample!


$$\text{Var}[\hat{\mathcal{L}}] = \frac{N^2}{M^2} \left(\sum_{i=1}^M \text{Var}[L_i] + 2 \sum_i \sum_i \text{Cov}[L_i, L_j] \right)$$
$$= N^2 \left(\frac{1}{M} \text{Var}[L_i] + \frac{M-1}{M} \text{Cov}[L_i, L_j] \right)$$

Kingma, Diederik P., Tim Salimans, and Max Welling. "Variational dropout and the local reparameterization trick." *Advances in NIPS* 2015.

The local reparameterization trick

Consider a linear layer with weight matrix W , input A and output B .

$$w_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma_{ij}^2)$$

$$B = AW$$

Predictions have **high**
correlation because there is
one weight sample per
batch

$$\mathbb{E}B = A\mu \quad \text{Var}B = A^2\sigma^2$$

$$B \sim \mathcal{N}(A\mu, A^2\sigma^2)$$

$$B = A\mu + \sqrt{A^2\sigma^2} \odot \epsilon$$

Predictions have **zero**
correlation because there is
one weight sample per
object

Kingma, Diederik P., Tim Salimans, and Max Welling. "Variational dropout and the local reparameterization trick." *Advances in NIPS* 2015.

The local reparameterization trick

LRT also reduces the variance of the stochastic gradient for **one** object

$\frac{\partial L}{\partial \mu_i}$ is the same for both RT and LRT, but

$$\frac{\partial L}{\partial \sigma_i} = \frac{\partial L}{\partial b} \cdot \frac{\partial b}{\partial \sigma_i} = \frac{\partial L}{\partial b} \cdot a_i \epsilon_i$$

RT, 1 sample per weight
A lot of redundant stochasticity

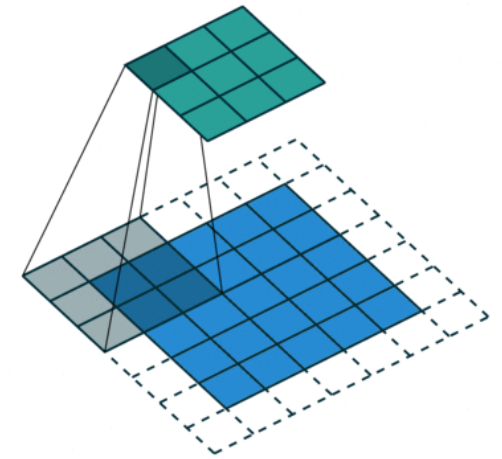
$$\frac{\partial L}{\partial \sigma_i} = \frac{\partial L}{\partial b} \cdot \frac{\partial b}{\partial \sigma_i} = \frac{\partial L}{\partial b} \cdot \frac{a_i^2 \sigma_i \epsilon}{\sqrt{a^2{}^\top \sigma^2}}$$

LRT, 1 sample per neuron
No redundant stochasticity

Kingma, Diederik P., Tim Salimans, and Max Welling. "Variational dropout and the local reparameterization trick." *Advances in NIPS* 2015.

LRT for convolutions

- B no longer factorizes in convolutional layers
 - Same weight samples should be used for different spatial positions
- Exact local reparameterization is too complex
 - We need to calculate the full covariance matrix for each activation
- We can use the mean-field local reparameterization as an approximation
 - Not justified (*yet*)
 - Performs much better than plain reparameterization



$$\begin{aligned}\mathbb{E}B &= A \star \mu & \text{Var}B &= A^2 \star \sigma^2 \\ B &\sim \mathcal{N}(A \star \mu, A^2 \star \sigma^2) \\ B &= A \star \mu + \sqrt{A^2 \star \sigma^2} \odot \epsilon\end{aligned}$$

What next?

- How to choose prior
- Faster test-time averaging
- Better posterior approximations
 - Implicit models

A good time take a 5 minute break?

Treating deterministic parameters

What about other parameters, not just weight matrices?

- Biases
- Linear transformation in BatchNorm
- Any other “non-expressive” parameters

1) We can put priors over them, and treat them as random variables

2) We can treat them as deterministic parameters

- Essentially it assumes a flat prior and a delta-peak posterior
- Or we could see it as bounding the marginal likelihood of the data given these parameters

$$\log p(\mathbf{t} | X, \boldsymbol{\theta}) \geq \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{w} | \boldsymbol{\phi})} \log p(\mathbf{t} | X, \mathbf{w}, \boldsymbol{\theta}) - \text{KL}(q(\mathbf{w} | \boldsymbol{\phi}) \parallel p(\mathbf{w})) \rightarrow \max_{\boldsymbol{\phi}, \boldsymbol{\theta}}$$

Empirical Bayes for Bayesian NNs

- How to choose the prior distribution?
- Type-II maximum likelihood (maximum evidence):

$$\begin{aligned} \log p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) &\rightarrow \max_{\boldsymbol{\theta}} \\ \log p(\mathbf{t}|\mathbf{X}, \boldsymbol{\theta}) &\geq \\ \geq \mathcal{L}(\boldsymbol{\phi}, \boldsymbol{\theta}) &= \mathbb{E}_{q(\mathbf{w}|\boldsymbol{\phi})} \log p(\mathbf{t}|\mathbf{X}, \mathbf{w}) - \text{KL}(q(\mathbf{w}|\boldsymbol{\phi})||p(\mathbf{w}|\boldsymbol{\theta})) \rightarrow \max_{\boldsymbol{\phi}, \boldsymbol{\theta}} \end{aligned}$$

- It is okay when $\dim \boldsymbol{\theta}$ is small
- May overfit if $\dim \boldsymbol{\theta}$ is large
 - Ideally we would have $p(\mathbf{w}|\boldsymbol{\theta}) = q(\mathbf{w}|\boldsymbol{\theta}) = \delta(\mathbf{w}_{ML})$
 - You never know whether you can overfit with a particular parameterization
- Usually used to induce sparsity (RVM, SWS, ...)

Distillation

Test-time averaging is expensive

Imaging we have a good sampler $q_t(w_t)$ for \mathbf{w}

- SG MCMC
- Variational approximate posterior

We can train a separate deterministic neural network (student) to “mimic” the ensemble (teacher):

$$\mathcal{L}(\mathbf{w}_{st}) = \mathbb{E}_{q_t(\mathbf{w}_t)} \mathcal{H}(p(\mathbf{t}|X, \mathbf{w}_{st}), p(\mathbf{t}|X, \mathbf{w}_t))$$

- Worse than the ensemble
- Better than a single network

Balan, Anoop Korattikara, et al. "Bayesian dark knowledge." *NIPS* 2015.

Distillation: examples

VI distillation:

- 1) Train the “teacher” network, obtain approx. posterior $q(\mathbf{w})$
- 2) For each iteration of learning of the “student” network:
 - 1) Sample a minibatch of data
 - 2) Sample predictions of the “teacher”
 - 3) Use SoftMax output of the teacher as “soft” labels

MCMC distillation:

- 1) Warm-up a Markov Chain for the “teacher”
- 2) For each iteration of learning of the “student” network:
 - 1) Sample a minibatch of data
 - 2) Make one SG MCMC update of the teacher
 - 3) Use SoftMax output of the teacher as “soft” labels

Variational inference with implicit posteriors

Implicit posterior: $\mathbf{w} = f(\epsilon, \phi)$, $q_\phi(\mathbf{w})$ is intractable

Example: $f(\epsilon, \phi)$ is an arbitrary neural network

$$\mathcal{L}(\phi) = \mathbb{E}_{q_\phi(\mathbf{w})} \log p(\mathbf{t}, \mathbf{w} | X) - \mathbb{E}_{q_\phi(\mathbf{w})} \log q_\phi(\mathbf{w})$$

Semi-implicit formulation:

$$q_\phi(\mathbf{w}) = \int q_\phi(\mathbf{w} | \mathbf{z}) q_\phi(\mathbf{z}) d\mathbf{z}$$

$$\mathbf{z} \sim q_\phi(\mathbf{z}), \quad \mathbf{w} \sim q_\phi(\mathbf{w} | \mathbf{z}) \Leftrightarrow \mathbf{w} \sim q_\phi(\mathbf{w})$$

- Any implicit distribution has a semi-implicit formulation $q(\mathbf{w} | \mathbf{z}) = \delta(\mathbf{w} - \mathbf{z})$
- Any semi-implicit distribution is implicit

Variational inference with implicit posteriors

	MCMC	IPM	Variational Bayes
Bias	No	???	Strong
Sampling/ Ensembling	Inefficient	???	Efficient
Density	No	???	Yes
Likelihood	Needed	???	Needed

Variational inference with implicit posteriors

	MCMC	IPM	Variational Bayes
Bias	No	Weak	Strong
Sampling/ Ensembling	Inefficient	Efficient	Efficient
Density	No	Can be estimated?	Yes
Likelihood	Needed	Can be avoided	Needed

Hierarchical variational inference

How to perform approximate inference with a semi-implicit posterior?

You should already know this:

- VI with auxiliary variables
- VI in RL with options

Assumptions:

- Can compute densities of both $q(w|z)$ and $q(z)$
- Can reparameterize both $q(w|z)$ and $q(z)$
- Can approximate inverse model $r(z|w) \approx \frac{q(w|z)q(z)}{q(w)}$

$$\mathcal{L} \geq \mathbb{E}_{q(w)} \log p(t, w|X) - \mathbb{E}_{q(z)q(w|z)} [\log q(w|z)q(z) - \log r(z|w)]$$

(Prove it if you didn't do it yet!)

Hierarchical variational inference

$$\mathcal{L}(\boldsymbol{\phi}) \geq \underline{\mathcal{L}}(\phi, \theta) = \mathbb{E}_{q_{\phi}(\mathbf{w})} \log p(\mathbf{t}, \mathbf{w} | X) - \mathbb{E}_{q_{\phi}(\mathbf{w}, \mathbf{z})} \log \frac{q_{\phi}(\mathbf{w} | \mathbf{z}) q_{\phi}(\mathbf{z})}{r_{\theta}(\mathbf{z} | \mathbf{w})} \rightarrow \max_{\phi, \theta}$$

Additional inference gap:


$$\mathcal{L}(\boldsymbol{\phi}) - \underline{\mathcal{L}}(\phi, \theta) = \mathbb{E}_{q_{\phi}(\mathbf{w})} \text{KL} \left(r_{\theta}(\mathbf{z} | \mathbf{w}) \parallel q_{\phi}(\mathbf{z} | \mathbf{w}) \right)$$

Multiplicative normalizing flows

Consider the following semi-implicit posterior:

$$\begin{aligned} q_{\phi}(\mathbf{w}) &= \int q_{\phi}(\mathbf{w}|\mathbf{z})q(\mathbf{z})d\mathbf{z} \\ q_{\phi}(w_{ij}|z_i) &= \mathcal{N}(w_{ij}|\mu_{ij}z_i, \sigma_{ij}^2) \\ q_{\phi}(\mathbf{z}) &= \text{NF}(\boldsymbol{\epsilon}, \boldsymbol{\phi}); \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, I) \\ r_{\theta}(z|w) &= \text{NF}(f_{\theta}(w), \theta) \end{aligned}$$

Can reparameterize and
compute log-density



- Learn non-trivial correlations between neurons
- Not clear how to evaluate the inference gap of HVI
- Normalizing flows are limited
 - Need a lot of depth for complex distributions
 - The width of each layer is the same
- One of the most complex posteriors for Bayesian NNs *yet*

Semi-implicit variational inference

- The rough idea: obtain a “MC” estimate of $q(w)$

$$\begin{aligned} q(w) &= \int q(w|z)q(z)dz = \mathbb{E}_{q(z)} q(w|z) \approx \\ &\approx q_{1..K}(w|z^1, \dots, z^K) = \frac{1}{K} \sum_{k=1}^K q(w|z^k) \end{aligned}$$

- Note that it is asymptotically exact:

$$q_{1..K}(w|z^1, \dots, z^K) \xrightarrow{K \rightarrow \infty} q(w)$$

- This idea can be used to obtain a sandwich bound for the entropy of $q(w)$:

$$\begin{aligned} -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{\textcolor{red}{K}} \sum_{k=1}^K q(w|z^k) &\geq \\ &\geq -\mathbb{E}_{q(w)} \log q(w) \geq \\ &\geq -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{\textcolor{red}{K} + 1} \sum_{k=0}^K q(w|z^k) \end{aligned}$$

SIVI: upper bound

$$-\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{K} \sum_{k=1}^K q(w|z^k) \geq -\mathbb{E}_{q(w)} \log q(w)$$

This is just plain Jensen's inequality ($\log \mathbb{E} \geq \mathbb{E} \log$)

$$\begin{aligned} -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{K} \sum_{k=1}^K q(w|z^k) &= -\mathbb{E}_{z^0, w|z^0} \mathbb{E}_{z^1 \dots z^K} \log \frac{1}{K} \sum_{k=1}^K q(w|z^k) \geq \\ &\geq -\mathbb{E}_{z^0, w|z^0} \log \frac{1}{K} \sum_{k=1}^K \mathbb{E}_{z^k} q(w|z^k) = -\mathbb{E}_{q(w)} \log q(w) \end{aligned}$$

SIVI: lower bound

$$-\mathbb{E}_{q(w)} \log q(w) \geq -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{K+1} \sum_{k=0}^K q(w|z^k)$$

Symmetrize the right side:

$$\begin{aligned} -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{K+1} \sum_{k=0}^K q(w|z^k) &= -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^j} \log \frac{1}{K+1} \sum_{k=0}^K q(w|z^k) = \\ &= -\frac{1}{K+1} \sum_{j=0}^K \mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{w|z^j} \log \frac{1}{K+1} \sum_{k=0}^K q(w|z^k) = \\ &= -\mathbb{E}_{z^0, \dots, z^K} \int \frac{1}{K+1} \sum_{j=0}^K q(w|z^j) \log \frac{1}{K+1} \sum_{k=0}^K q(w|z^k) dw = \\ &= -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log q_{0..K}(w|z^0, \dots, z^K) \end{aligned}$$

SIVI: lower bound

$$-\mathbb{E}_{q(w)} \log q(w) \geq -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log q_{0..K}(w|z^0, \dots, z^K)$$

Rewrite the left side in the same expectations:

$$-\mathbb{E}_{q(w)} \log q(w) = -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log q(w)$$

And subtract the right side:

$$\begin{aligned} & -\mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log q(w) + \\ & + \mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log q_{0..K}(w|z^0, \dots, z^K) = \\ & = \mathbb{E}_{z^0, \dots, z^K} \mathbb{E}_{q_{0..K}(w|z^0, \dots, z^K)} \log \frac{q_{0..K}(w|z^0, \dots, z^K)}{q(w)} = \\ & = \mathbb{E}_{z^0, \dots, z^K} \text{KL}(q_{0..K}(w|z^0, \dots, z^K) \parallel q(w)) \geq 0 \end{aligned}$$

Semi-implicit variational inference

- This idea can be used to obtain a sandwich for ELBO

- $\mathcal{L}_K \leq \mathcal{L} \leq \mathcal{L}^K$
- Both \mathcal{L}_K and \mathcal{L}^K monotonically converge to \mathcal{L}
- We can now estimate ELBO in **any** semi-implicit model!
- $q(z)$ can be fully implicit (**any** reparameterizable distribution)

Optimize w.r.t. ϕ

$$\underline{\mathcal{L}}_K = \mathbb{E}_w \log p(t, w|X) - \mathbb{E}_{z^0, z^1, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{\underline{K} + 1} \sum_{k=0}^K q(w|z^k)$$

$$\overline{\mathcal{L}}^K = \mathbb{E}_w \log p(t, w|X) - \mathbb{E}_{z^0, z^1, \dots, z^K} \mathbb{E}_{w|z^0} \log \frac{1}{\overline{K}} \sum_{k=1}^K q(w|z^k)$$

Bound the
inference gap

Implicit models: takeaways

- We can now model arbitrarily complex posteriors!
- Applicable to both Bayesian NNs and VAEs
- The next big thing in Bayesian learning!



Bayesian neural networks: takeaways

- Regularization by noise
- Reparameterization
- Local reparameterization
- Deterministic parameters
- Empirical Bayes
- Use implicit model if you want better approximation

Unbiased implicit variational inference

Main idea: we can consider the reparameterization of $q_\phi(z)$ as a part of the conditional $q_\phi(w|z)$:

$$\int q_\phi(w|z)q_\phi(z)dz = \int q_\phi(w|\epsilon)p(\epsilon)d\epsilon$$

$$q_\phi(w|\epsilon) = q_\phi(w|z) \Big|_{z=f_\phi(\epsilon)}$$

The model stays the same, but the joint distribution is now explicit!

$$q_\phi(w, \epsilon) = q_\phi(w|\epsilon)p(\epsilon)$$

We now might be able to run HVI with a fully implicit $q_\phi(z)$...

... Or go even further

Unbiased implicit variational inference

Provides an unbiased gradient estimate for the ELBO

$$\begin{aligned}\nabla_{\phi} \mathbb{E}_{q_{\phi}(w)} \log q_{\phi}(w) &= \nabla_{\phi} \mathbb{E}_{p(\epsilon)} \log q_{\phi}(f_{\phi}(\epsilon)) = \\ &= \mathbb{E}_{p(\epsilon)} \nabla_{\phi} \log q_{\phi}(w) \Big|_{w=f_{\phi}(\epsilon)} + \mathbb{E}_{p(\epsilon)} \nabla_w \log q_{\phi}(w) \Big|_{w=f_{\phi}(\epsilon)} \cdot \nabla_{\phi} f_{\phi}(\epsilon) = \\ &= \underbrace{\mathbb{E}_{q_{\phi}(w)} \nabla_{\phi} \log q_{\phi}(w)}_{=0} + \mathbb{E}_{p(\epsilon)} \nabla_w \log q_{\phi}(w) \Big|_{w=f_{\phi}(\epsilon)} \cdot \nabla_{\phi} f_{\phi}(\epsilon) \\ \nabla_w \log q_{\phi}(w) &= \mathbb{E}_{q(\epsilon'|w)} \nabla_{\phi} \log q_{\phi}(w|\epsilon')\end{aligned}$$

Can be sampled using plain HMC!

We can start HMC using pair (ϵ, w) to avoid warm-up

Unbiased implicit variational inference: proof?

$$\nabla_w \log q_\phi(w) = \mathbb{E}_{q(\epsilon'|w)} \nabla_w \log q_\phi(w|\epsilon')$$

$$\begin{aligned} \nabla_w \log q_\phi(w) &= \frac{1}{q_\phi(w)} \nabla_w \int q_\phi(w|\epsilon') q(\epsilon') d\epsilon' = \frac{1}{q_\phi(w)} \int \nabla_w q_\phi(w|\epsilon') q(\epsilon') d\epsilon' = \\ &= \frac{1}{q_\phi(w)} \int q_\phi(w|\epsilon') q(\epsilon') \nabla_w \log q_\phi(w|\epsilon') d\epsilon' = \\ &= \int \frac{q_\phi(w|\epsilon') q(\epsilon')}{q_\phi(w)} \nabla_w \log q_\phi(w|\epsilon') d\epsilon' = \\ &= \int q(\epsilon'|w) \nabla_w \log q_\phi(w|\epsilon') d\epsilon' \end{aligned}$$

Unbiased implicit variational inference

Provides an unbiased gradient estimate for the ELBO

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(w)} \log q_{\phi}(w) = \mathbb{E}_{p(\epsilon)} \nabla_w \log q_{\phi}(w) \Big|_{w=f_{\phi}(\epsilon)} \cdot \nabla_{\phi} f_{\phi}(\epsilon)$$
$$\nabla_w \log q_{\phi}(w) = \mathbb{E}_{q(\epsilon|w)} \nabla_w \log q_{\phi}(w|\epsilon)$$

- Still cannot estimate ELBO directly (only gradients)
- Can use SIVI bounds to monitor / compare ELBO
- Need sufficiently “nice” $q(\epsilon)$
 - $\mathcal{N}(0, I)$ works fine and is used almost always
 - Discrete $q(\epsilon)$ may be problematic (e.g. no dropout in the implicit generator f_{ϕ})