Bayesian Inference and Deep Learning

Dr. Richard E. Turner (ret26@cam.ac.uk) Computational and Biological Learning Lab, Department of Engineering, University of Cambridge

Limitations of Deep Learning

Robust Deep Learning

fragile (adversarial examples)

well calibrated uncertainty estimates: deep learning is often confidently wrong

Data-Efficient Deep Learning

small data, big models (few-shot learning and reinforcement learning) leverage heterogenous data sources (multi-task learning)

Flexible Deep Learning

continual learning (online learning & model building) active learning (RL exploration-exploitation)

correct diff ostrich 0.0 0.2 0.4 0.6 0.8 1.0 prediction confidence

frequency



ResNet (2016) CIFAR-100



Logistic Regression as a motivating example

Logistic regression: A single neuron



Logistic regression: Maximum Likelihood Estimation

observe data: estimate parameters

data: estimate parameters

$$\mathcal{X} = \{\mathbf{x}_n\}_{n=1}^N$$

$$\mathcal{Y} = \{y_n\}_{n=1}^N$$

$$\overset{\circ}{=} 0$$

$$\overset{\circ}{=} 1$$

$$\frac{\mathbf{x} \times \mathbf{x} \times \mathbf{x} \times \mathbf{x}} = \mathbf{x}_0 + \mathbf{x}_1$$

$$\mathbf{x} = \mathbf{x}_0 + \mathbf{x}_1$$

maximum likelihood estimate: parameters that make observed data most probable

$$P(\mathcal{Y}|\mathbf{w}, \mathcal{X}) = \prod_{n=1}^{N} P(y_n | \mathbf{w}, \mathbf{x}_n) = \prod_{n=1}^{N} f(\mathbf{w}^{\top} \mathbf{x}_n)^{y_n} (1 - f(\mathbf{w}^{\top} \mathbf{x}_n))^{(1-y_n)}$$
$$\mathbf{w}^{\mathsf{ML}} = \underset{\mathbf{w}}{\operatorname{arg\,max}} \log P(\mathcal{Y}|\mathbf{w}, \mathcal{X}) = \underset{\mathbf{w}}{\operatorname{arg\,max}} \sum_{n=1}^{N} \left[y_n \log f(\mathbf{w}^{\top} \mathbf{x}_n) + (1 - y_n) \log(1 - f(\mathbf{w}^{\top} \mathbf{x}_n)) \right]$$

Logistic regression: Maximum Likelihood learning

observe data: estimate parameters

data and prediction objective $\log P(\mathcal{Y}|\mathbf{w}, \mathcal{X})$ 0 0.0 1.0 $P(y = 1 | \mathbf{w}, \mathbf{x})$ $\log P(\mathcal{Y}|\mathbf{w}, \mathcal{X})$ -2.5 -50 0 0.5 o =-5.0 w_1 1 -7.5 × 0.0 -10.050 100 -2 4 Ó -10 -5 Ó 0 2 iteration w_0 \mathbf{x}

model:

 $\mathbf{w}^{\mathsf{ML}} = \underset{\mathbf{w}}{\arg\max} \log P(\mathcal{Y}|\mathbf{w}, \mathcal{X}) = \arg\max_{\mathbf{w}} \sum_{n=1}^{N} \left[y_n \log f(\mathbf{w}^{\top} \mathbf{x}_n) + (1 - y_n) \log(1 - f(\mathbf{w}^{\top} \mathbf{x}_n)) \right]$ 6/31



observe data: estimate parameters



$$\begin{aligned} & \text{model:} \\ f(\mathbf{x}, \mathbf{w}) = P\left(y = 1 | \mathbf{w}, \mathbf{x}\right) = \frac{1}{1 + \exp(-\mathbf{w}^{\top}\mathbf{x})} = \frac{1}{1 + \exp(-(w_0 + w_1 x_1))} \\ & \text{maximum likelihood estimate: parameters that make observed data most probable} \end{aligned}$$

$$\mathbf{w}^{\mathsf{ML}} = \underset{\mathbf{w}}{\operatorname{arg\,max}} \log P(\mathcal{Y}|\mathbf{w}, \mathcal{X}) = \underset{\mathbf{w}}{\operatorname{arg\,max}} \sum_{n=1}^{N} \left[y_n \log f(\mathbf{w}^{\top} \mathbf{x}_n) + (1 - y_n) \log(1 - f(\mathbf{w}^{\top} \mathbf{x}_n)) \right]_{6/31}$$

Bayesian approaches to logistic regression

Probabilistic model

encodes prior assumptions in a recipe for generating datasets

$$p(\mathbf{w}) = \mathcal{G}(\mathbf{w}; \mathbf{0}, \sigma_{\mathbf{w}}^2 \mathbf{I})$$

$$P(y=1|\mathbf{w},\mathbf{x}) = \frac{1}{1 + \exp(-a(\mathbf{x},\mathbf{w}))}$$

Probabilistic inference

1. "right" answer is a probability distribution over all possible settings of the weights that indicates the plausiblity of that setting given data



only way to be coherent, Cox 1946

only way to protect against Dutch books, Ramsey 1926

2. apply the sum and product rules of probability to compute the plausibility of any setting of any unknown variable



The posterior distribution over weights



The predictive distribution over class labels







11/31

Bayesian Inference in Action: 1D Classification Example





Bayesian Inference in Action: 1D Classification Example

Bayesian Inference in Action: 1D Classification Example





Bayesian Inference in Action: 1D Classification Example

12/31

Why be Bayesian (distributional estimates of weights)?

	maximum-likelihood	Bayes	
learning	$\mathbf{w}^{\scriptscriptstyleML} = rg\max_{\mathbf{w}} \log P(\mathcal{Y} \mathbf{w}, \mathcal{X})$	$p(\mathbf{w} \mathcal{Y}, \mathcal{X}) = \frac{1}{P(\mathcal{Y} \mathcal{X})} p(\mathbf{w}) P(\mathcal{Y} \mathbf{w}, \mathcal{X})$	
prediction	$p(y^* x^*,\mathcal{Y},\mathcal{X}) \approx p(y^* \mathbf{w}^{\mathrm{ML}},x^*)$	$p(y^* x^*, \mathcal{Y}, \mathcal{X}) = \int p(\mathbf{w} \mathcal{Y}, \mathcal{X}) p(y^* \mathbf{w}, x^*) d\mathbf{w}$	w
	single weight setting	ensemble over requir weight settings approximation	es on
Robust D	eep Learning		

point estimates over-confident, averaging over weight settings less so Bayesian methods are more robust to adversarial examples (hard to fool ensemble of networks + uncertainty)

Data-efficient Deep Learning

small data, big model: build models 'the size of a house' & let data prune/learn structure leverage heterogeneous data sources (multi-task learning) using shared parameters

Flexibile Deep Learning

continual learning: use old posterior as prior active learning: select data that are expected to reduce uncertainty in parameter estimates the most

Approximate Bayesian Inference

Laplace's approximation: MacKay 1991 (Saddle point approximation)

$$p(\mathbf{w}|\mathcal{Y},\mathcal{X}) \propto p(\mathcal{Y},\mathbf{w}|\mathcal{X}) = p^{*}(\mathbf{w})$$

$$\mathbf{w}^{\text{MAP}} = \underset{\mathbf{w}}{\text{arg max }} \log p(\mathbf{w}|\mathcal{Y},\mathcal{X}) = \underset{\mathbf{w}}{\text{arg max }} \left[\log p(\mathcal{Y}|\mathbf{w},\mathcal{X}) + \log p(\mathbf{w})\right]$$
Taylor expand log-prob to 2nd order about MAP est:
$$\log p^{*}(\mathbf{w}) \approx \log p^{*}(\mathbf{w}^{\text{MAP}}) + \underbrace{\mathbf{w}^{\text{MAP}} + \frac{1}{2}(\mathbf{w} - \mathbf{w}^{\text{MAP}})^{\top} \frac{d^{2}\log p^{*}(\mathbf{w})}{d\mathbf{w}d\mathbf{w}^{\top}}\Big|_{\mathbf{w} = \mathbf{w}^{\text{MAP}}} \mathbf{w}$$

$$p^{*}(\mathbf{w}) \approx p^{*}(\mathbf{w}^{\text{MAP}}) \exp\left(\frac{1}{2}(\mathbf{w} - \mathbf{w}^{\text{MAP}})^{\top} \frac{d^{2}\log p^{*}(\mathbf{w})}{d\mathbf{w}d\mathbf{w}^{\top}}\Big|_{\mathbf{w} = \mathbf{w}^{\text{MAP}}} (\mathbf{w} - \mathbf{w}^{\text{MAP}})\right)$$

$$p^{*}(\mathbf{w}) \approx p^{*}(\mathbf{w}^{\text{MAP}}) \det(2\pi\Sigma)^{1/2}\mathcal{G}(\mathbf{w}; \mathbf{w}^{\text{MAP}}, \Sigma)$$

$$\Sigma^{-1} = -\frac{d^{2}\log p^{*}(\mathbf{w})}{d\mathbf{w}d\mathbf{w}^{\top}}\Big|_{\mathbf{w} = \mathbf{w}^{\text{MAP}}} \mathbf{w}^{(\mathbf{w})} \cdot q(\mathbf{w})$$
Prediction requires additional approx.: $p(y^{*}|\mathbf{x}^{*}, \mathcal{X}, \mathcal{Y}) = \int p(y^{*}|\mathbf{w}, \mathbf{x}^{*})q(\mathbf{w})d\mathbf{w} \approx \frac{1}{M}\sum_{m=1}^{M} p(y^{*}|\mathbf{w}^{(m)}, \mathbf{x}^{*})$

$$\frac{15/31}{2}$$

Laplace's approximation: MacKay 1991 (Saddle point approximation)



Laplace's approximation: Classification Example



Variational Inference: the KL Divergence

Kullback-Leibler (KL) divergence

1. non-negative

$$\frac{\delta^2}{\delta q^2} \mathrm{KL}(q(\mathbf{w}) || p(\mathbf{w})) = \frac{1}{q(\mathbf{w})} \ge 0$$

2. zero (minimised) when $q(\mathbf{w}) = p(\mathbf{w})$

$$\mathrm{KL}(p(\mathbf{w})||p(\mathbf{w})) = \int p(\mathbf{w}) \log \frac{p(\mathbf{w})}{p(\mathbf{w})} \mathrm{d}\mathbf{w} = 0$$



Variational Inference: the KL Divergence

Kullback–Leibler (KL) divergence

1. non-negative

_

$$\frac{\delta^2}{\delta q^2} \mathrm{KL}(q(\mathbf{w}) || p(\mathbf{w})) = \frac{1}{q(\mathbf{w})} \ge 0$$

2. zero (minimised) when $q(\mathbf{w}) = p(\mathbf{w})$ $\mathrm{KL}(p(\mathbf{w})||p(\mathbf{w})) = \int p(\mathbf{w}) \log \frac{p(\mathbf{w})}{p(\mathbf{w})} \mathrm{d}\mathbf{w} = 0$

divergence measures 'distance' between distributions

3. can be computed up to an additive constant w/o needing normalisation for $p(\mathbf{w})$

$$\operatorname{KL}\left(q(\mathbf{w})||\frac{1}{Z}p^{*}(\mathbf{w})\right) = \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{\frac{1}{Z}p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} = \int q(\mathbf{w})\log\frac{q(\mathbf{w})Z}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w}$$
$$= \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} + \int q(\mathbf{w})\log Z\mathrm{d}\mathbf{w} = \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} + \log Z$$
$$\longrightarrow \text{ suitable for approximate inference}$$

 $\mathrm{KL}(q(\mathbf{w})||p(\mathbf{w})) = \int q(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w})} \mathrm{d}\mathbf{w}$

Variational Inference

Kullback-Leibler (KL) divergence $KL(q(\mathbf{w})||p(\mathbf{w})) = \int q(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w})} d\mathbf{w} \stackrel{1\&2}{\geq 0}$ equality when
 $q(\mathbf{w}) = p(\mathbf{w})$

Use KL to optimise an approximation $q(\mathbf{w})$ to the true posterior

$$\begin{split} q^{\mathrm{VI}}(\mathbf{w}) &= \mathop{\mathrm{arg\,min}}_{q \in \mathcal{Q}} \operatorname{KL}(q(\mathbf{w}) || p(\mathbf{w} | \mathcal{Y}, \mathcal{X})) \stackrel{\mathbf{3}}{=} \mathop{\mathrm{arg\,min}}_{q \in \mathcal{Q}} \int q(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w}) p(\mathcal{Y} | \mathbf{w}, \mathcal{X})} \mathrm{d}\mathbf{w} \\ &\approx \mathop{\mathrm{arg\,min}}_{q \in \mathcal{Q}} \frac{1}{M} \sum_{m=1}^{M} \log \frac{q(\mathbf{w}^{(m)})}{p(\mathbf{w}^{(m)}) p(\mathcal{Y} | \mathbf{w}^{(m)}, \mathcal{X})} \quad \mathbf{w}^{(m)} \sim q(\mathbf{w}) \xrightarrow{\operatorname{requires}}_{\substack{\text{reparameterisation}\\ \operatorname{trick}}} \end{split}$$

3. can be computed up to an additive constant w/o needing normalisation for $p(\mathbf{w})$

$$\operatorname{KL}\left(q(\mathbf{w})||\frac{1}{Z}p^{*}(\mathbf{w})\right) = \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{\frac{1}{Z}p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} = \int q(\mathbf{w})\log\frac{q(\mathbf{w})Z}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w}$$
$$= \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} + \int q(\mathbf{w})\log Z\mathrm{d}\mathbf{w} = \int q(\mathbf{w})\log\frac{q(\mathbf{w})}{p^{*}(\mathbf{w})}\mathrm{d}\mathbf{w} + \log Z$$
$$\longrightarrow \text{ suitable for approximate inference}$$

19/31

Variational Inference: Classification Example

$$q^{\text{VI}}(\mathbf{w}) = \underset{q \in \mathcal{Q}}{\operatorname{arg\,min}} \operatorname{KL}(q(\mathbf{w}) || p(\mathbf{w} | \mathcal{Y}, \mathcal{X}))$$

$$\approx \underset{q \in \mathcal{Q}}{\operatorname{arg\,min}} \frac{1}{M} \sum_{m=1}^{M} \log \frac{q(\mathbf{w}^{(m)})}{p(\mathbf{w}^{(m)})p(\mathcal{Y} | \mathbf{w}^{(m)}, \mathcal{X})}$$

$$\underset{factorised\,Gaussian}{\operatorname{approximate\,family:}} W_{1}$$

$$q(\mathbf{w}) = \mathcal{G}(w_{0}; \mu_{0}, \sigma_{0}^{2}) \times \mathcal{G}(w_{1}; \mu_{1}, \sigma_{1}^{2})$$

$$\underset{\mu_{0}}{\operatorname{optimise\,w.r.t.}} \mu_{0}, \sigma_{0}^{2}, \mu_{1}, \sigma_{1}^{2}}$$

VI converts inference into optimization

20/31

 σ_1^2

 w_0



Variational Inference vs Laplace: Classification Example

Case study 1. Robust Deep Learning Uncertainty calibration & adversarial examples





Case study 2. Flexible Deep Learning Continual multi-task learning



25 / 31

What is Continual Learning?





A zoo of discriminative continual learning tasks

Continual Learning Test 1: Permuted MNIST (online non-iid inputs, single head)





Summary

- **Continual learning** is naturally handled by **Bayesian inference**: allows multi-task transfer and avoids catastrophic forgetting
- Variational Continual Learning is a state-of-the-art continual learning method
- Orthogonal research directions: **complex models** (adapting more than just the head of the network) and **online automatic model building**

Variational Continual Learning, ICLR 2018 Streaming Sparse Gaussian Process Approximations, NIPS 2017

Case study 3: Data Efficient Deep Learning One-shot learning

30/31

One shot learning using Approximate Bayesian inference



One shot learning using Approximate Bayesian inference



One shot learning using Approximate Bayesian inference

