The connection of dropout and Bayesian statistics

Interpretation of dropout as approximate Bayesian modelling of NN

http://mlg.eng.cam.ac.uk/yarin/thesis/thesis.pdf



Uncertainty in Deep Learning



Yarin Gal

Department of Engineering University of Cambridge

October 2016

This dissertation is submitted for the degree of $Doctor \ of \ Philosophy$



Geoffrey Hinton Google, University Toronto

Dropout



Srivastava et al., *Journal of Machine Learning Research* 15 (2014)



Nitish Srivastava University Toronto

Presented by Beate Sick, IDP, ZHAW

Outline

- Main principles of frequentists and Bayesians a reminder
- Neural networks and how they get trained a reminder
- What is dropout and what is the connection to Bayesian inference?
- How to do approximate Bayesian inference in deep NN
- How to use dropouts to get an uncertainty measure for predictions?
- How can we use the Bayesian approach to make better models?





Frequentist's and Bayesian's view on data and model parameters

Frequentist:

- Data are a repeatable random sample
 -> e.g. 10 coin tosses: frequencies 2xtail, 8xheads Model: Bernoulli with head-probability as parameter θ
- Underlying model parameters are fixed during this repeatable data sampling
- We use Maximum Likelihood to determine the parameter value under which the observed data is most likely.
- 95% CI means that when repeating the experiment many times 95% of all 95% CI will cover the true fixed parameter.

Bayesian:

- Data are observed once and seen as fixed
- Parameters are described probabilistically
- Before data observation we formulate a prior belief in parameter distribution
- After data observation we updated posterior distribution of model parameters.
- 95% credibility interval means an interval of parameter values that cover 95% of the posterior distribution.

Distribution of data under fixed parameter



Bayesian modeling has less problems with complex models

Frequentist's strategy:

You can only use a complex model if you have enough data!

Bayesian's strategy:

Use the model complexity you believe in.

Do not just use the best fitting model.

Do use the full posterior distribution over parameter settings leading to vague predictions since many different parameter settings have significant posterior probability.

 $p(\mathbf{y} | \mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y} | \mathbf{x}, \boldsymbol{\omega}) \cdot p(\boldsymbol{\omega} | \mathbf{X}, \mathbf{Y}) d\boldsymbol{\omega}$

prediction via marginalization over ω :



Logistic regression or a neural net with 1 neuron



A neural network is a parametric model

We observe

- inputs
$$\mathbf{X} = {\{\mathbf{x}_i\}}_{i=1,...,N}$$

- outputs
$$\mathbf{Y} = \{y_i\}_{i=1,...,N}$$

Using a NN as output generating model we can do parametric inference for the connection weights W in the NN that are the parameter defining the model.

The output is $y = f(x, W) = W_2(\sigma(W_1x + b))$



Fitting a NN classifier via softmax likelihood optimization



Update weights in back-pass



chain rule (via backpropagation)

$$P(y=k) = p_k = \frac{e^{z_k}}{\sum_{j \in \text{class.id}} e^{z_j}}$$

LogLikelihood = cost-function C
 Find sets of weights θ that minimize C!

$$C=-\sum_k I_k \log \, p_k\,$$
 , ${
m I_k}\,$ is indicator for true class k



- Main principles of frequentists and Bayesians a reminder
- Neural networks and how they get trained a reminder
- What is dropout and what is the connection to Bayesian inference?
- How to do approximate Bayesian inference in deep NN
- How to use dropouts to get an uncertainty measure for predictions?
- How can we use the Bayesian approach to make better models?





Dropout



"dropout"

At each training step we remove random nodes with a probability of p resulting in a sparse version of the full net and we use backpropagation to update the weights.

-> In each training step we train another NN model, but all models share weights!

Meaning a certain connection weight is the same in all models at the current update value.

Why "dropout" can be a good idea

Dropout reduces complexity of the model and thus overfitting ?



Using dropout during training implies:

- We train in each training step another sparse model. Only weights to not-dropped units are updated
- Sharing weights may introduce "regularization"
- By averaging over these models we should be able to "reduce noise", "overfitting"



Use the trained net without dropout during test time

But we need to do a small adjustment!

Why "dropout" can be a good idea



Use the trained net without dropout during test time

Q: Suppose that with all inputs present at test time and the output of this neuron is x.

What would its output be during training time, in expectation? (e.g. if p = 0.5)

during test: $a = w0^*x + w1^*y$

during training $E[a] = \frac{1}{4} * (w0*0 + w1*0 + w0*0 + w1*0 + w0*0 + w1*y + w0*x + w1*0 + w0*x + w1*x + w0*x + w1*x + w0*x + w$



Using all units in the test time forward pass leads to an output that 2x the expected output during training (this is true for linear units without hidden layers).

=> Have to compensate by reducing the weights during test time by multiplying them by the dropout probability p=0.5

 $w0^{*}x + w1^{*}y)$

 $= \frac{1}{4} * (2 \text{ w0} * x + 2 \text{ w1} * y)$

 $= \frac{1}{2} * (w0*x + w1*y)$

Why "dropout" can be a good idea

The training data consists of many different pictures of Oliver Dürr and Albert Einstein



We need a huge number of neurons to extract good features which help to distinguish Oliver from Einstein



Dropout forces the network to learn redundant and independent features



We will see that dropout is an approximation of full Bayesian learning using a Bayesian Network



Zurich University of Applied Sciences





"Bayesian NN" At each training step we update the posterior distribution of the weights In test time we use the posterior to determine the predictive distribution

"dropout" At each training step we remove random nodes with a probability p

Bayesian Modeling in the framework of NN

- In Bayesian Modeling we define a prior distribution over the parameter W: $W_i \sim N(0, I)$ defining $p(w_i)$
- For classification NN tasks we assume a softmax likelihood

$$P(\mathbf{y}=\mathbf{k} \mid \boldsymbol{\omega}, \mathbf{x}) = \mathbf{p}_{k} = \frac{e^{\mathbf{f}_{\mathbf{k},\boldsymbol{\omega}}(\mathbf{x})}}{\sum_{k'} e^{\mathbf{f}_{\mathbf{k},\boldsymbol{\omega}}(\mathbf{x})}}$$



• Given a dataset X, Y we then look for the posteriori distribution capturing the most probable model parameter given the observed data

likelihood prior

$$p(\omega | \mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y} | \omega, \mathbf{X}) \cdot p(\omega)}{p(\mathbf{Y} | \mathbf{X})}$$

normalizer=marginal likelihood

Marginalization steps in Bayesian statistics

• We can use the posterior distribution

 $p(\mathbf{w}|\mathbf{X},\mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{w},\mathbf{X}) \cdot p(\mathbf{w})}{p(\mathbf{Y}|\mathbf{X})}$

- -> to determine the most probable weight vector MAP: Maximum a Posteriori
- -> to determine credible intervals
- -> sample weight vectors from posterior distribution



 For the posteriori we need to determine the normalizer, also called model evidence via integrating over all possible parameters weighted by their probabilities. This process is also called marginalization over ω leading to the name marginal likelihood

 $p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X},\omega) \cdot p(\omega) d\omega$

Using a neural network for prediction

A deep NN or deep CNN

- Performs a hierarchical non-linear transformation of its input y = f(x, W)
- These models give us only point estimates with no uncertainty information.

Remark: for shallow NN CIs for the weights can be determined if the network is identified (White, 1998)





With Bayesian modeling we can get a measure of uncertainty by evaluating the posterior distribution of the NN weights.

- Main principles of frequentists and Bayesians a reminder
- Neural networks and how they get trained a reminder
- What is dropout and what is the connection to Bayesian inference?
- How to do approximate Bayesian inference in deep NN
- How to use dropouts to get an uncertainty measure for predictions?
- How can we use the Bayesian approach to make better models?





Full Bayesian learning and ideas for approximations

Full Bayesian learning means computing the full posterior distribution over all possible parameter settings.

- This is extremely computationally intensive for all but the simplest models
- But it allows us to use complicated models with not so much data

Approximation of full Bayesian learning is the only way if we have models with many (>100) parameters or complex models.

- We can use Markov Chain Monte Carlo (MCMC) methods to sample parameter vectors ω according to their posterior probabilities at the observed data $p(\omega|X,Y)$ and use them to estimate the distribution $q(\omega)$.
- We can approximate the posterior distribution for the model parameters via

a) Integrated Nested Laplace Approximations (INLA)

- b) Variational inference
 - replacing the posterior distribution at the observed data $p(\omega|X, Y)$ with a member $q(\omega)$ of a simpler distribution family Q that minimizes the Kullback–Leibler divergence to the posterior

Variational inference approximates full Bayesian learning

- Approximate $p(\omega | \mathbf{X}; \mathbf{Y})$ with simple distribution $q_{\theta}(\omega)$
- Minimize Kullback Leibler divergence of q from the posterior w.r.t. to the variational parameters θ:

$$\operatorname{KL}(q_{\theta}(\omega) \| \mathbf{p}(\omega | \mathbf{X}, \mathbf{Y})) = \int q_{\theta}(\omega) \log \frac{q_{\theta}(\omega)}{\mathbf{p}(\omega | \mathbf{X}, \mathbf{Y})} d\omega = E_q \left[\log \left(q_{\theta}(\omega) \right) - \log \left(\mathbf{p}(\omega | \mathbf{X}, \mathbf{Y}) \right) \right]$$

Minimizing KL means to wiggle the parameter θ of q to find the value of θ for which q resembles the posterior distribution as good as possible.



The term variational is used because you pick the **best** q in Q -- the term derives from the "calculus of variations," which deals with optimization problems. A particular q in Q is specified by setting some variational parameters -- the parameter which is adjusted during optimization.

Variational inference to approximate Bayesian learning ctd.

Minimizing the KL divergence of q from the posterior dist. p w.r.t. θ

$$\mathrm{KL}(q_{\theta}(\omega) \| \mathbf{p}(\omega | \mathbf{X}, \mathbf{Y})) = \int q_{\theta}(\omega) \log \frac{q_{\theta}(\omega)}{\mathbf{p}(\omega | \mathbf{X}, \mathbf{Y})} d\omega$$

is equivalent to maximizing a lower bound of the log marginal likelihood w.r.t. θ (also called *evidence lower bound* or ELBO)

$$L_{VI}(\theta) \coloneqq \int q_{\theta}(\omega) \cdot \log \left(p(\mathbf{Y}|\mathbf{X}, \omega) \right) d\omega - KL(q_{\theta}(\omega) \| \mathbf{p}(\omega)) \leq L = \log \left(p(\mathbf{Y} | \mathbf{X}) \right)$$

The key trick is here like in most variational methods, to write the true log-likelihood L as the log of an expectation under some q. We can then get a lower bound via Jensen's inequality, which tells us that log expectation >= expectation log (since log is concave and q is a distribution).

$$\begin{split} \mathbf{L} &= \log\left(p(\mathbf{Y} \mid \mathbf{X})\right) = \log\int p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right) \cdot p\left(\omega\right) d\omega = \log\int p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right) \cdot p\left(\omega\right) \frac{q_{\theta}(\omega)}{q_{\theta}(\omega)} d\omega = \log\left(E_{q_{\theta}}\left[\frac{p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right) \cdot p\left(\omega\right)}{q_{\theta}(\omega)}\right]\right] \\ &\geq E_{q_{\theta}}\left[\log\left(\frac{p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right) \cdot p\left(\omega\right)}{q_{\theta}(\omega)}\right)\right] = E_{q_{\theta}}\left[\log\left(p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right)\right) + \log\left(\frac{p\left(\omega\right)}{q_{\theta}(\omega)}\right)\right] = E_{q_{\theta}}\left[\log\left(p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right)\right)\right] - E_{q_{\theta}}\left[\log\left(\frac{q_{\theta}(\omega)}{p\left(\omega\right)}\right)\right] \\ &= \int q_{\theta}(\omega) \cdot \log\left(p\left(\mathbf{Y} \mid \mathbf{X}, \omega\right)\right) d\omega - KL\left(q_{\theta}(\omega) \mid p(\omega)\right) \end{split}$$

Variational inference to approximate Bayesian learning ctd.

The best θ that makes q_{θ} a good approximation for the posterior is given by the θ that maximizes the following objective as lower bound of L:

$$L_{VI}(\theta) \coloneqq \int q_{\theta}(\omega) \cdot \log \left(p(\mathbf{Y}|\mathbf{X}, \omega) \right) d\omega - KL(q_{\theta}(\omega) \| \mathbf{p}(\omega))$$

Since this integral is not tractable for almost all q therefore we will MC integration to approximate this quantity.

We will sample $\hat{\omega}$ from q and look at the following term where for each sampling step the integral is replaced by $\log(p(Y|X, \hat{\omega}))$ leading to a new objective:

$$\hat{\mathbf{L}}(\theta) \coloneqq \log\left(p\left(\mathbf{Y}|\mathbf{X},\hat{\omega}\right)\right) - KL\left(q_{\theta}(\omega) \| \mathbf{p}(\omega)\right)$$

In approximate Bayesian statistics it is known, that \hat{L} is an unbiased estimator of L.

Variational inference to approximate Bayesian learning ctd.

$$\begin{split} \mathbf{L}_{VI}(\theta) &\coloneqq \int q_{\theta}(\omega) \cdot \log \Big(p \big(\mathbf{Y} | \mathbf{X}, \omega \big) \Big) d\omega - KL \big(q_{\theta}(\omega) \, \| \, \mathbf{p}(\omega) \big) \\ \hat{\mathbf{L}}(\theta) &\coloneqq \log \Big(p \big(\mathbf{Y} | \mathbf{X}, \hat{\omega} \big) \Big) - KL \big(q_{\theta}(\omega) \, \| \, \mathbf{p}(\omega) \big) \end{split}$$

Results from stochastic inference guarantee that optimizing \hat{L} w.r.t. θ will converge to the same optimal θ than optimizing L_{VI} w.r.t. θ .

To stepwise optimize \hat{L} we sample in each step one $\hat{\omega}$ from $q_{\theta}(\omega)$ and do one step of optimization of θ and then we update $q_{\theta}(\omega)$ before sampling the next $\hat{\omega}$.

For inference repeatedly do:

- a) Sample $\widehat{\omega} \sim q_{\theta}(\omega)$
- b) Do one step of minimization w.r.t. θ :

 $\hat{\mathbf{L}}(\theta) \coloneqq \log\left(p\left(\mathbf{Y}|\mathbf{X}, \hat{\omega}\right)\right) - KL\left(q_{\theta}(\omega) \| \mathbf{p}(\omega)\right)$

What kind of q-distribution should we use to resemble dropout with Bayesian learning?

Define the structure of the approximate distribution q

We define q to factorizes over weight matrices \mathbf{W}_{i} of the layers 1 to L.

$$q_{\theta}(\omega) = \prod_{i} q_{\mathbf{M}_{i}}(\mathbf{W}_{i})$$

For each \mathbf{W}_i we define q to be the product of the mean weight matrix \mathbf{M}_i and a diagonal matrix which holds Bernoulli variables on the diagonal.





M_i is as variational parameter of q

Sampling the diagonal elements **z** from a Bernoulli is identical to randomly setting columns of M to zero which is identical to randomly setting units of the network to zero -> dropout!

Why this structure of the approximate distribution q

- it resembles dropout!
 We know what to do during test time: use full NN and set W_i = p_iW_i
- Bernoullis are computationally cheap to get multi-modality
- This q constrains the weights to be near zero:
 - for this q structure the variance of the weights are:

$$\operatorname{Var}(\mathbf{W}_i) = \mathbf{M}_i \cdot \mathbf{M}_i^T \cdot \mathbf{p}_i \cdot (1 - \mathbf{p}_i)$$

 we know that posterior uncertainty decreases with more data implying that for fixed dropout probability p_i the average weights need to decrease with more data.

The strongest regularization is achieved with dropout probability $p_i = 0.5$.

- Main principles of frequentists and Bayesians a reminder
- Neural networks and how they get trained a reminder
- What is dropout and what is the connection to Bayesian inference?
- How to do approximate Bayesian inference in deep NN
- How to use dropouts to get an uncertainty measure for predictions?
- How can we use the Bayesian approach to make better models?





DL with dropout as approximate Bayesian inference allows to get a prediction distribution for uncertainty estimation

To get an approximation of the posterior via training

- Randomly set columns of M_i to zero (do dropout)
- Update the weights by doing one step



To sample from the learned approximate posterior we just can do dropout during the test time when using the trained NN for prediction.

From the received predictions we can estimate the predictive distribution and from this different uncertainty measures such as the variance.

Can this insights help us to make better models?

Yes – just do dropout also in the convolutional layer during training and test time!

MC dropout is equivalent to performing *several* stochastic forward passes through the network and averaging the results. By doing so Yarin Gal was able to outperform state of the art error rates on MNIST and CIFAR-10 without changing the architecture of the used CNNs or anything else beside dropout.



Dropout approach speeds up reinforcement learning



Example: train a Roomba vacuum cleaner via enforcement learning

Without dropout: take random action with probability ϵ and the model-predicted action otherwise (Epsilon-greedy)

With dropout: Do several stochastic forward passes through the dropout network and choose action with highest value (Thompson sampling)

Thank you for your attention



How good is the dropout approach?

Compare performance of dropout-method with well established and published variational inference (VI) and probabilistic back-propagation (PBP) results on the basis of root-mean-square-error (RMSE) and log-likelihood (LL).

The better fitting the model is the smaler is the RMSE.

	Avg. Test RMSE and Std. Errors			Avg. Test LL and Std. Errors		
Dataset	VĬ	PBP	Dropout	VI	PBP	Dropout
Boston Housing	4.32 ± 0.29	3.01 ±0.18	2.97 ± 0.85	-2.90 ± 0.07	-2.57 ±0.09	-2.46 ±0.25
Concrete Strength	7.19 ± 0.12	5.67 ± 0.09	$\textbf{5.23} \pm \textbf{0.53}$	-3.39 ± 0.02	-3.16 ±0.02	-3.04 ±0.09
Energy Efficiency	2.65 ± 0.08	1.80 ± 0.05	1.66 ± 0.19	-2.39 ± 0.03	-2.04 ± 0.02	-1.99 ±0.09
Kin8nm	0.10 ± 0.00	$\textbf{0.10} \pm \textbf{0.00}$	$\textbf{0.10} \pm \textbf{0.00}$	0.90 ± 0.01	0. 9 0 ±0. 01	0.95 ±0.03
Naval Propulsion	0.01 ± 0.00	$\textbf{0.01} \pm \textbf{0.00}$	$\textbf{0.01} \pm \textbf{0.00}$	3.73 ± 0.12	3.73 ± 0.01	3.80 ± 0.05
Power Plant	4.33 ± 0.04	4.12 ± 0.03	$\textbf{4.02} \pm \textbf{0.18}$	-2.89 ± 0.01	-2.84 ± 0.01	-2.80 ± 0.05
Protein Structure	4.84 ± 0.03	4.73 ±0.01	$\textbf{4.36} \pm \textbf{0.04}$	-2.99 ± 0.01	-2.97 ± 0.00	-2.89 ± 0.01
Wine Quality Red	0.65 ± 0.01	0.64 ± 0 .01	$\textbf{0.62} \pm \textbf{0.04}$	-0.98 ± 0.01	-0.97 ±0.01	-0.93 ±0.06
Yacht Hydrodynamics	6.89 ± 0.67	1.02 ± 0.05	1.11 ± 0.38	-3.43 ± 0.16	-1.63 ± 0.02	-1.55 ±0.12
Year Prediction MSD	$9.034 \pm NA$	$8.879 \pm NA$	$8.849 \pm NA$	$-3.622 \pm NA$	$-3.603 \pm NA$	-3.588 ±NA

The smaller the uncertainty of the model is the larger is the LL.

Table 1: Average test performance in RMSE and predictive log likelihood for a popular variational inference method (VI, Graves [20]), Probabilistic back-propagation (PBP, Hernández-Lobato and Adams [27]), and dropout uncertainty (Dropout).