

# Dropout in RNNs Following a VI Interpretation

Yarin Gal

[yg279@cam.ac.uk](mailto:yg279@cam.ac.uk)

Recurrent neural networks (RNNs) are damn useful.

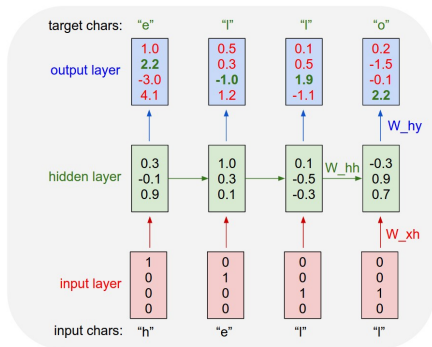


Figure : RNN structure

Image Source:

[karpathy.github.io/2015/05/21/rnn-effectiveness](https://karpathy.github.io/2015/05/21/rnn-effectiveness)

But these also overfit very quickly...

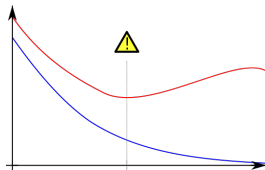


Figure : Overfitting

This means...

- ▶ We can't use large models
- ▶ We have to use early stopping
- ▶ We can't use small data
- ▶ We have to waste data for validation sets...

But these also overfit very quickly...

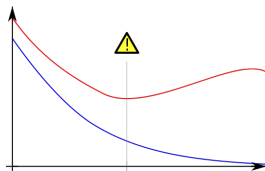


Figure : Overfitting

This means...

- ▶ We can't use large models
- ▶ We have to use early stopping
- ▶ We can't use small data
- ▶ We have to waste data for validation sets...

But these also overfit very quickly...

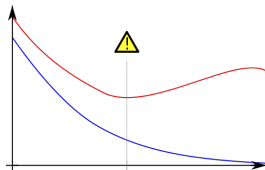


Figure : Overfitting

This means...

- ▶ We can't use large models
- ▶ We have to use early stopping
- ▶ We can't use small data
- ▶ We have to waste data for validation sets...

But these also overfit very quickly...

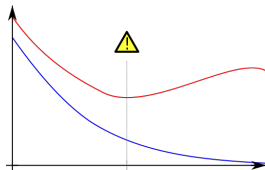


Figure : Overfitting

This means...

- ▶ We can't use large models
- ▶ We have to use early stopping
- ▶ We can't use small data
- ▶ We have to waste data for validation sets...

But these also overfit very quickly...

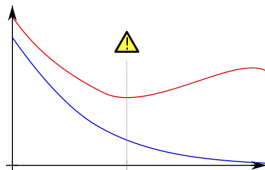


Figure : Overfitting

This means...

- ▶ We can't use large models
- ▶ We have to use early stopping
- ▶ We can't use small data
- ▶ We have to waste data for validation sets...

Let's use dropout then. But lots of research has claimed that that's a **bad idea**:

- ▶ **Pachitariu & Sahani, 2013**
  - ▶ noise added in the recurrent connections of an RNN leads to model instabilities
- ▶ **Bayer et al., 2013**
  - ▶ with dropout, the RNNs dynamics change dramatically
- ▶ **Pham et al., 2014**
  - ▶ dropout in recurrent layers disrupts the RNNs ability to model sequences
- ▶ **Zaremba et al., 2014**
  - ▶ applying dropout to the non-recurrent connections alone results in improved performance
- ▶ **Bluche et al., 2015**
  - ▶ exploratory analysis of the performance of dropout before, inside, and after the RNNs



→ has settled on using dropout for **inputs and outputs alone**:

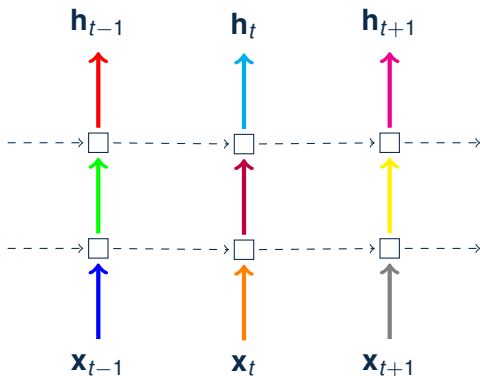


Figure : Naive application of dropout in RNNs (colours = different dropout masks)

## Why not use dropout with recurrent layers?

- ▶ ~~It doesn't work~~
- ▶ ~~Noise drowns the signal~~
- ▶ Because it's not used correctly?

Why not use dropout with recurrent layers?

- ▶ ~~It doesn't work~~
- ▶ ~~Noise drowns the signal~~
- ▶ Because it's not used correctly?

Why not use dropout with recurrent layers?

- ▶ ~~It doesn't work~~
- ▶ ~~Noise drowns the signal~~
- ▶ Because it's not used correctly?

Why not use dropout with recurrent layers?

- ▶ ~~It doesn't work~~
- ▶ ~~Noise drowns the signal~~
- ▶ Because it's not used correctly?

**First, some background on Bayesian modelling and VI in Bayesian neural networks.**

- ▶ Observed inputs  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and outputs  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N$
- ▶ Capture stochastic process believed to have generated outputs
- ▶ Def.  $\omega$  model parameters as r.v.
- ▶ Prior dist. over  $\omega$ :  $p(\omega)$
- ▶ Likelihood:  $p(\mathbf{Y}|\omega, \mathbf{X})$
- ▶ Posterior:  $p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\omega, \mathbf{X})p(\omega)}{p(\mathbf{Y}|\mathbf{X})}$  (Bayes' theorem)
- ▶ Predictive distribution given new input  $\mathbf{x}^*$

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) \underbrace{p(\omega|\mathbf{X}, \mathbf{Y})}_{\text{posterior}} d\omega$$

- ▶ But...  $p(\omega|\mathbf{X}, \mathbf{Y})$  is often intractable

- ▶ Observed inputs  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and outputs  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N$
- ▶ Capture stochastic process believed to have generated outputs
- ▶ Def.  $\omega$  model parameters as r.v.
- ▶ Prior dist. over  $\omega$ :  $p(\omega)$
- ▶ Likelihood:  $p(\mathbf{Y}|\omega, \mathbf{X})$
- ▶ Posterior:  $p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\omega, \mathbf{X})p(\omega)}{p(\mathbf{Y}|\mathbf{X})}$  (Bayes' theorem)
- ▶ Predictive distribution given new input  $\mathbf{x}^*$

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) \underbrace{p(\omega|\mathbf{X}, \mathbf{Y})}_{\text{posterior}} d\omega$$

- ▶ But...  $p(\omega|\mathbf{X}, \mathbf{Y})$  is often intractable

- ▶ Observed inputs  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N$  and outputs  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N$
- ▶ Capture stochastic process believed to have generated outputs
- ▶ Def.  $\omega$  model parameters as r.v.
- ▶ Prior dist. over  $\omega$ :  $p(\omega)$
- ▶ Likelihood:  $p(\mathbf{Y}|\omega, \mathbf{X})$
- ▶ Posterior:  $p(\omega|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\omega, \mathbf{X})p(\omega)}{p(\mathbf{Y}|\mathbf{X})}$  (Bayes' theorem)
- ▶ Predictive distribution given new input  $\mathbf{x}^*$

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) \underbrace{p(\omega|\mathbf{X}, \mathbf{Y})}_{\text{posterior}} d\omega$$

- ▶ But...  $p(\omega|\mathbf{X}, \mathbf{Y})$  is often intractable



- ▶ Approximate  $p(\omega|\mathbf{X}, \mathbf{Y})$  with simple dist.  $q_\theta(\omega)$
- ▶ Minimise divergence from posterior w.r.t.  $\theta$

$$\text{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

- ▶ Identical to minimising

$$\mathcal{L}_{\text{VI}}(\theta) := - \int q_\theta(\omega) \log \overbrace{p(\mathbf{Y}|\mathbf{X}, \omega)}^{\text{likelihood}} d\omega + \text{KL}(q_\theta(\omega) \parallel \overbrace{p(\omega)}^{\text{prior}})$$

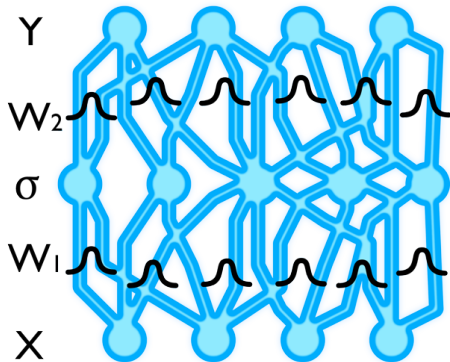
- ▶ We can approximate the **predictive distribution**

$$q_\theta(\mathbf{y}^*|\mathbf{x}^*) = \int p(\mathbf{y}^*|\mathbf{x}^*, \omega) q_\theta(\omega) d\omega.$$

- ▶ Place prior  $p(\mathbf{W}_i)$ :

$$\mathbf{W}_i \sim \mathcal{N}(0, \mathbf{I})$$

for  $i \leq L$  (and write  $\omega := \{\mathbf{W}_i\}_{i=1}^L$ ).



- ▶ Output is a r.v.  $\mathbf{f}(\mathbf{x}, \omega) = \mathbf{W}_L \sigma(\dots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots)$ .

- ▶ Place prior  $p(\mathbf{W}_i)$ :

$$\mathbf{W}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

for  $i \leq L$  (and write  $\omega := \{\mathbf{W}_i\}_{i=1}^L$ ).

- ▶ Output is a r.v.  $\mathbf{f}(\mathbf{x}, \omega) = \mathbf{W}_L \sigma(\dots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots)$ .
- ▶ Softmax likelihood for class.:  $p(y|\mathbf{x}, \omega) = \text{softmax}(\mathbf{f}(\mathbf{x}, \omega))$   
or a Gaussian for regression:  $p(\mathbf{y}|\mathbf{x}, \omega) = \mathcal{N}(\mathbf{y}; \mathbf{f}(\mathbf{x}, \omega), \tau^{-1}\mathbf{I})$ .
- ▶ But difficult to evaluate posterior

$$p(\omega|\mathbf{X}, \mathbf{Y}).$$

- ▶ Place prior  $p(\mathbf{W}_i)$ :

$$\mathbf{W}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

for  $i \leq L$  (and write  $\omega := \{\mathbf{W}_i\}_{i=1}^L$ ).

- ▶ Output is a r.v.  $\mathbf{f}(\mathbf{x}, \omega) = \mathbf{W}_L \sigma(\dots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots)$ .
- ▶ Softmax likelihood for class.:  $p(y|\mathbf{x}, \omega) = \text{softmax}(\mathbf{f}(\mathbf{x}, \omega))$   
or a Gaussian for regression:  $p(\mathbf{y}|\mathbf{x}, \omega) = \mathcal{N}(\mathbf{y}; \mathbf{f}(\mathbf{x}, \omega), \tau^{-1} \mathbf{I})$ .
- ▶ But difficult to evaluate posterior

$$p(\omega|\mathbf{X}, \mathbf{Y}).$$

- ▶ Place prior  $p(\mathbf{W}_i)$ :

$$\mathbf{W}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

for  $i \leq L$  (and write  $\omega := \{\mathbf{W}_i\}_{i=1}^L$ ).

- ▶ Output is a r.v.  $\mathbf{f}(\mathbf{x}, \omega) = \mathbf{W}_L \sigma(\dots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \dots)$ .
- ▶ Softmax likelihood for class.:  $p(y|\mathbf{x}, \omega) = \text{softmax}(\mathbf{f}(\mathbf{x}, \omega))$   
or a Gaussian for regression:  $p(\mathbf{y}|\mathbf{x}, \omega) = \mathcal{N}(\mathbf{y}; \mathbf{f}(\mathbf{x}, \omega), \tau^{-1} \mathbf{I})$ .
- ▶ But difficult to evaluate posterior

$$p(\omega|\mathbf{X}, \mathbf{Y}).$$

- ▶ Def  $q_\theta(\omega)$  to approximate posterior  $p(\omega|\mathbf{X}, \mathbf{Y})$
- ▶ KL divergence to minimise:

$$\text{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

$$\propto \boxed{-\int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega} + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$

$$=: \mathcal{L}(\theta)$$

- ▶ Approximate the integral with MC integration  $\hat{\omega} \sim q_\theta(\omega)$ :

$$\hat{\mathcal{L}}(\theta) := -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$

- ▶ Def  $q_\theta(\omega)$  to approximate posterior  $p(\omega|\mathbf{X}, \mathbf{Y})$
- ▶ KL divergence to minimise:

$$\text{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

$$\propto \boxed{-\int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega} + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$

$$=: \mathcal{L}(\theta)$$

- ▶ Approximate the integral with MC integration  $\hat{\omega} \sim q_\theta(\omega)$ :

$$\hat{\mathcal{L}}(\theta) := -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$

- ▶ Def  $q_\theta(\omega)$  to approximate posterior  $p(\omega|\mathbf{X}, \mathbf{Y})$
- ▶ KL divergence to minimise:

$$\text{KL}(q_\theta(\omega) \parallel p(\omega|\mathbf{X}, \mathbf{Y}))$$

$$\propto \boxed{-\int q_\theta(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega} + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$

$$=: \mathcal{L}(\theta)$$

- ▶ Approximate the integral with MC integration  $\hat{\omega} \sim q_\theta(\omega)$ :

$$\hat{\mathcal{L}}(\theta) := -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_\theta(\omega) \parallel p(\omega))$$



- ▶ Unbiased estimator:

$$E_{\hat{\omega} \sim q_{\theta}(\omega)}(\hat{\mathcal{L}}(\theta)) = \mathcal{L}(\theta)$$

- ▶ Converges to the same optima as  $\mathcal{L}(\theta)$
- ▶ For inference, repeat:
  - ▶ Sample  $\hat{\omega} \sim q_{\theta}(\omega)$
  - ▶ And minimise (one step)

$$\hat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta$ .

- ▶ Unbiased estimator:

$$E_{\hat{\omega} \sim q_{\theta}(\omega)}(\hat{\mathcal{L}}(\theta)) = \mathcal{L}(\theta)$$

- ▶ Converges to the same optima as  $\mathcal{L}(\theta)$
- ▶ For inference, repeat:
  - ▶ Sample  $\hat{\omega} \sim q_{\theta}(\omega)$
  - ▶ And minimise (one step)

$$\hat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta$ .

- ▶ Unbiased estimator:

$$E_{\hat{\omega} \sim q_{\theta}(\omega)}(\hat{\mathcal{L}}(\theta)) = \mathcal{L}(\theta)$$

- ▶ Converges to the same optima as  $\mathcal{L}(\theta)$
- ▶ For inference, repeat:
  - ▶ Sample  $\hat{\omega} \sim q_{\theta}(\omega)$
  - ▶ And minimise (one step)

$$\hat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta$ .

- ▶ Given variational parameters  $\theta = \{[\mathbf{m}_{i1}, \dots, \mathbf{m}_{iK}]\}_{i=1}^L$ :

$$q_{\theta}(\boldsymbol{\omega}) = \prod_i q_{\theta}(\mathbf{W}_i)$$

$$q_{\theta}(\mathbf{W}_i) = \prod_k q_{\mathbf{m}_{ik}}(\mathbf{w}_{ik})$$

$$q_{\mathbf{m}_{ik}}(\mathbf{w}_{ik}) = p\mathcal{N}(0, \sigma^2) + (1 - p)\mathcal{N}(\mathbf{m}_{ik}, \sigma^2)$$

→  $k$ 'th column of the  $i$ 'th layer is a multivariate mixture of Gaussians

- ▶ With small enough  $\sigma^2$ , in practice equivalent to

$$\mathbf{z}_{i,j} \sim \text{Bernoulli}(p_i) \text{ for } i = 1, \dots, L, j = 1, \dots, K_{i-1}$$

$$\mathbf{W}_i = \mathbf{M}_i \cdot \text{diag}([\mathbf{z}_{i,j}]_{j=1}^{K_i})$$

with  $\mathbf{z}_{i,j}$  Bernoulli r.v.s.

In summary:

Minimise divergence between  $q_{\theta}(\omega)$  and  $p(\omega|\mathbf{X}, \mathbf{Y})$ :

► Repeat:

► Sample  $\hat{\mathbf{z}}_{i,j} \sim \text{Bernoulli}(p_j)$  and set

$$\hat{\mathbf{W}}_i = \mathbf{M}_i \cdot \text{diag}([\hat{\mathbf{z}}_{i,j}]_{j=1}^{K_i})$$

$$\hat{\omega} = \{\hat{\mathbf{W}}_i\}_{i=1}^L$$

► Minimise (one step)

$$\hat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta = \{\mathbf{M}_i\}_{i=1}^L$  (set of matrices).

In summary:

Minimise divergence between  $q_{\theta}(\omega)$  and  $p(\omega|\mathbf{X}, \mathbf{Y})$ :

- ▶ Repeat:
  - ▶ = Randomly set columns of  $\mathbf{M}_i$  to zero
  - ▶ Minimise (one step)

$$\widehat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \widehat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta = \{\mathbf{M}_i\}_{i=1}^L$  (set of matrices).

In summary:

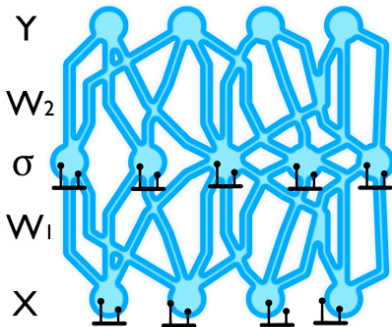
Minimise divergence between  $q_{\theta}(\omega)$  and  $p(\omega|\mathbf{X}, \mathbf{Y})$ :

- ▶ Repeat:
  - ▶ = Randomly set units of the network to zero
  - ▶ Minimise (one step)

$$\widehat{\mathcal{L}}(\theta) = -\log p(\mathbf{Y}|\mathbf{X}, \widehat{\omega}) + \text{KL}(q_{\theta}(\omega) \parallel p(\omega))$$

w.r.t.  $\theta = \{\mathbf{M}_i\}_{i=1}^L$  (set of matrices).

Sounds familiar?



$$\hat{\mathcal{L}}(\theta) = \underbrace{-\log p(\mathbf{Y}|\mathbf{X}, \hat{\omega})}_{= \text{loss}} + \underbrace{\text{KL}(q_{\theta}(\omega) \parallel p(\omega))}_{= L_2 \text{ reg}}$$

**Implementing VI with  $q_{\theta}(\cdot)$  above = implementing dropout in deep network**



- ▶ Multiplicative Gaussian noise (Srivastava et al. 2014) –
- ▶ Multiply network units by  $\mathcal{N}(1, 1)$
- ▶ Same performance as dropout



## Multiplicative Gaussian noise as approximate inference<sup>1</sup>

$$\mathbf{z}_{i,j} \sim \mathcal{N}(1, 1) \text{ for } i = 1, \dots, L, j = 1, \dots, K_{i-1}$$

$$\mathbf{W}_i = \mathbf{M}_i \cdot \text{diag}([\mathbf{z}_{i,j}]_{j=1}^{K_i})$$

$$q_{\theta}(\omega) = \prod q_{\mathbf{M}_i}(\mathbf{W}_i)$$

Similarly for **drop-connect** (Wan et al., 2013), etc.

---

<sup>1</sup>See Gal and Ghahramani (2015) and Kingma et al. (2015)

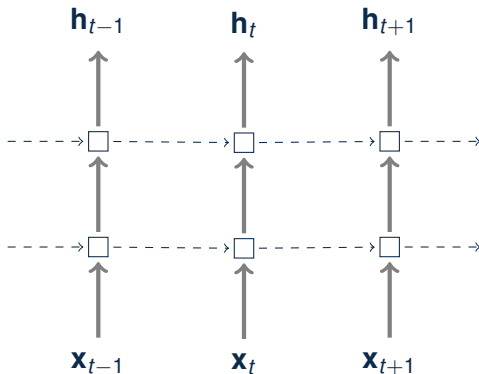


Figure : A Recurrent Neural Network

- ▶ Input sequence of vectors  $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$  with  $T$  time steps
- ▶ Let  $\omega = \{\text{all weight matrices in the model}\}$
- ▶ Define  $\mathbf{h}_t = \mathbf{f}_h^\omega(\mathbf{x}_t, \mathbf{h}_{t-1})$ 
  - ▶ **single** recurrent unit transition. E.g. tanh of affine transformation:  $\tanh(W\mathbf{x}_t + U\mathbf{h}_{t-1} + b)$
- ▶ Set  $\mathbf{f}_y^\omega(\mathbf{h}_T) = \mathbf{f}_y^\omega(\mathbf{f}_h^\omega(\mathbf{x}_T, \dots \mathbf{f}_h^\omega(\mathbf{x}_1, \mathbf{h}_0) \dots))$ 
  - ▶ model output (e.g. affine transformation of last state, or function of all states)
- ▶ Lastly, define  $p(\mathbf{y} | \mathbf{f}_y^\omega(\mathbf{h}_T))$ 
  - ▶ model likelihood. E.g.  $\mathcal{N}(\mathbf{y}; \mathbf{f}_y^\omega(\mathbf{h}_T), \sigma^2)$
- ▶ Similarly for LSTM, GRU

- ▶ Looking at the variational lower bound, we have:

$$\int q(\omega) \log p(\mathbf{y} | \mathbf{f}_y^\omega(\mathbf{h}_T)) d\omega =$$
$$\int q(\omega) \log p\left(\mathbf{y} \mid \mathbf{f}_y^\omega\left(\mathbf{f}_h^\omega(\mathbf{x}_T, \dots, \mathbf{f}_h^\omega(\mathbf{x}_1, \mathbf{h}_0) \dots)\right)\right) d\omega,$$

- ▶ Using MC integration with  $\hat{\omega} \sim q(\omega)$ ,

$$\mathcal{L}_{VI} \approx -\log p\left(\mathbf{y} \mid \mathbf{f}_y^{\hat{\omega}}\left(\mathbf{f}_h^{\hat{\omega}}(\mathbf{x}_T, \dots, \mathbf{f}_h^{\hat{\omega}}(\mathbf{x}_1, \mathbf{h}_0) \dots)\right)\right)$$
$$+ \text{KL}(q_\theta(\omega) \parallel p(\omega)).$$

Objective:

$$-\log p\left(\mathbf{y} \mid \mathbf{f}_y^{\hat{\omega}} \left( \mathbf{f}_h^{\hat{\omega}} \left( \mathbf{x}_T, \dots, \mathbf{f}_h^{\hat{\omega}} \left( \mathbf{x}_1, \mathbf{h}_0 \right) \dots \right) \right) \right) + \dots \quad \hat{\omega} \sim q(\omega)$$

- In practice, **use the same dropout mask at each time step**

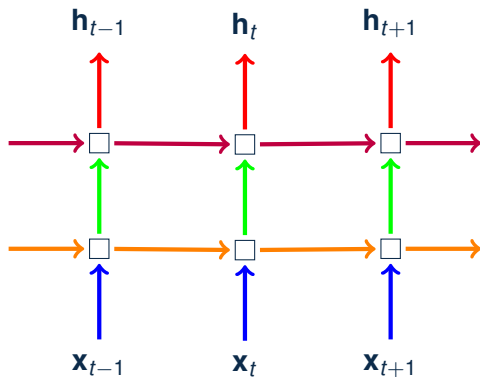


Figure : Bayesian motivated dropout in RNNs (colours = dropout masks)

- ▶ With **continuous** inputs we apply dropout to the input layer (place a distr. over weight matrix)
- ▶ But not for models with **discrete** inputs...
- ▶ Word embeddings: input can be seen as either the word embed. itself, or a “**one-hot**” encoding times an embed. matrix
- ▶ Optimising embedding matrix can lead to **overfitting**...
- ▶ Let's apply dropout **to the one-hot encoded vectors**

- ▶ **In practice**, drop words at random throughout the sentence
  - ▶ Randomly set embedding matrix **rows** to zero – entire word embeddings
  - ▶ Mask is **repeated** at each time step → **drop the same words throughout the sequence**
  - ▶ i.e. drop word **types** at random rather than word **tokens**
- ▶ For **example**, “the dog and the cat” might become “— dog and — cat” or “the — and the cat”, but never “— dog and the cat”.
- ▶ Can be interpreted as encouraging model to **not “depend” on single words.**

Some results (much more in paper):

- Sentiment analysis (Pang & Lee, 2005)

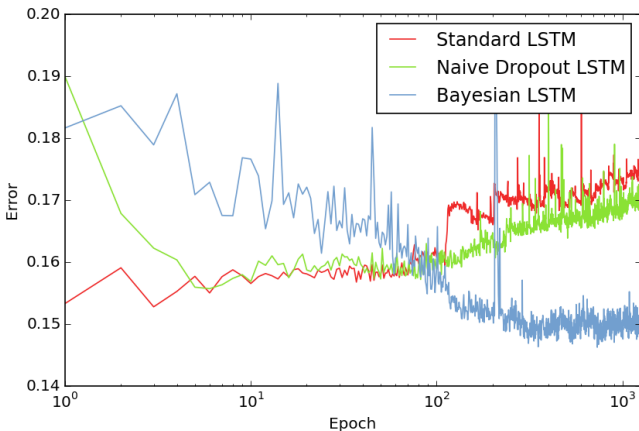


Figure : LSTM test error



Some results (much more in paper):

- Sentiment analysis (Pang & Lee, 2005)

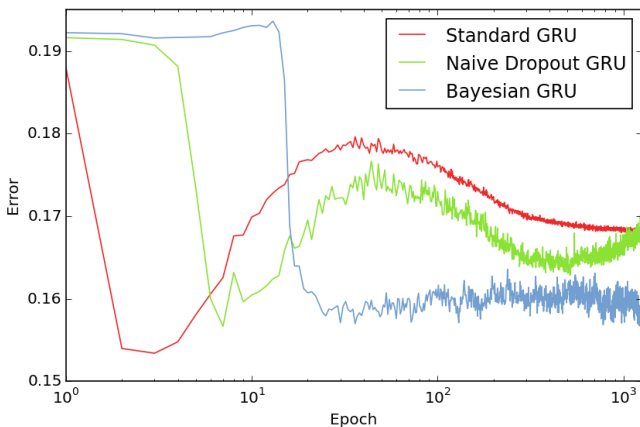


Figure : GRU test error

Some results (much more in paper):

- ▶ Sentiment analysis (Pang & Lee, 2005)
- ▶ Language model (Penn Treebank)

	Medium LSTM			Large LSTM		
	Validation	Test	WPS	Validation	Test	WPS
Non-regularized (early stopping)	121.1	121.7	5.5K	128.3	127.4	2.5K
Moon et al. [19]	100.7	97.0	4.8K	122.9	118.7	3K
Moon et al. [19] +emb dropout	88.9	86.5	4.8K	88.8	86.0	3K
Zaremba et al. [4]	86.2	82.7	5.5K	82.2	78.4	2.5K
Variational (tied weights)	81.8 ± 0.2	79.7 ± 0.1	4.7K	77.3 ± 0.2	75.0 ± 0.1	2.4K
Variational (tied weights, MC)	–	79.0 ± 0.1	–	–	74.1 ± 0.0	–
Variational (untied weights)	81.9 ± 0.2	79.7 ± 0.1	2.7K	77.9 ± 0.3	75.2 ± 0.2	1.6K
Variational (untied weights, MC)	–	<b>78.6 ± 0.1</b>	–	–	<b>73.4 ± 0.0</b>	–

Some results (much more in paper):

- ▶ Sentiment analysis (Pang & Lee, 2005)
- ▶ Language model (Penn Treebank)

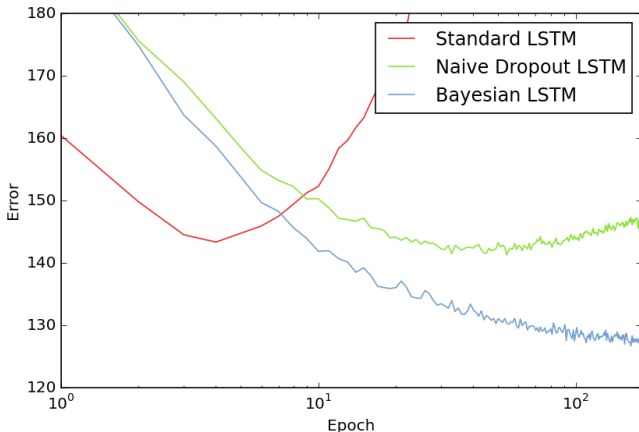


Figure : 2 layers LSTM, 200 units

- ▶ Practical deep learning uncertainty?
  - ▶ Capture language ambiguity?

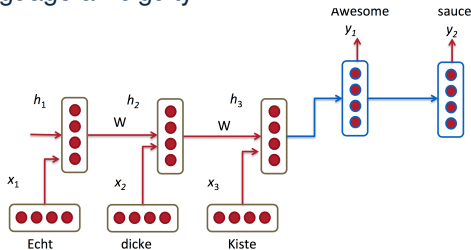
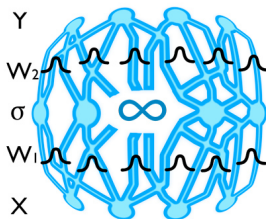


Image Source: [cs224d.stanford.edu/lectures/CS224d-Lecture8.pdf](https://cs224d.stanford.edu/lectures/CS224d-Lecture8.pdf)

- ▶ Weight uncertainty for model debugging?
- ▶ Principled extensions of deep learning?
  - ▶ New appr. distributions = new stochastic reg. techniques?
  - ▶ Model compression:  $W_i \sim$  discrete distribution w. continuous base measure?

- ▶ Practical deep learning uncertainty?
  - ▶ Capture language ambiguity?
  - ▶ Weight uncertainty for model debugging?



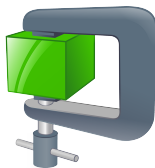
- ▶ Principled extensions of deep learning?
  - ▶ New appr. distributions = new stochastic reg. techniques?
  - ▶ Model compression:  $W_j \sim$  discrete distribution w. continuous base measure?

- ▶ Practical deep learning uncertainty?
  - ▶ Capture language ambiguity?
  - ▶ Weight uncertainty for model debugging?
- ▶ Principled extensions of deep learning?
  - ▶ New appr. distributions = new stochastic reg. techniques?

$$q_{\theta}(\omega) = ?$$

- ▶ Model compression:  $\mathbf{W}_i \sim$  discrete distribution w. continuous base measure?

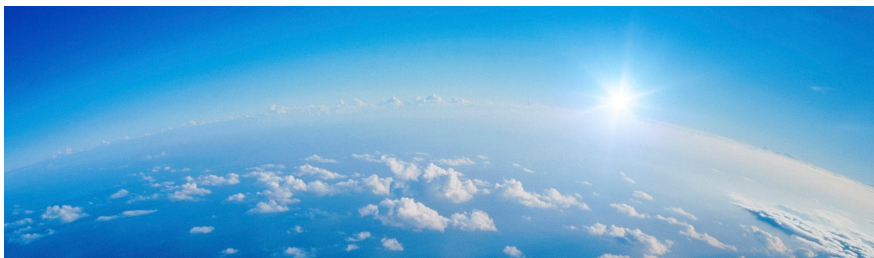
- ▶ Practical deep learning uncertainty?
  - ▶ Capture language ambiguity?
  - ▶ Weight uncertainty for model debugging?
- ▶ Principled extensions of deep learning?
  - ▶ New appr. distributions = new stochastic reg. techniques?
  - ▶ Model compression:  $\mathbf{W}_i \sim$  discrete distribution w. continuous base measure?



- ▶ Practical deep learning uncertainty?
  - ▶ Capture language ambiguity?
  - ▶ Weight uncertainty for model debugging?
- ▶ Principled extensions of deep learning?
  - ▶ New appr. distributions = new stochastic reg. techniques?
  - ▶ Model compression:  $\mathbf{W}_i \sim$  discrete distribution w. continuous base measure?

Work in progress!

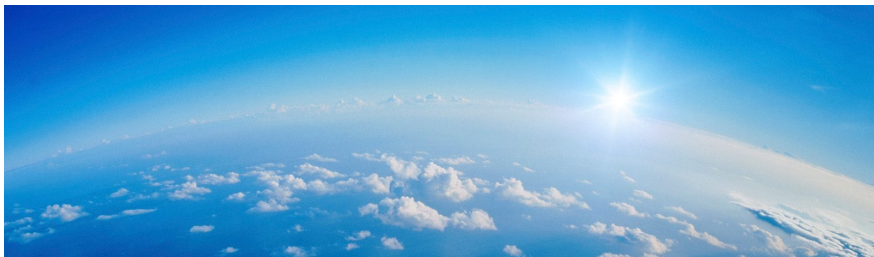




*Most exciting is work to come:*

- ▶ **Practical uncertainty** in deep learning applications
- ▶ **Principled extensions** to deep learning tools
- ▶ **Hybrid** deep learning – Bayesian models

*and much, much, more.*



*Most exciting is work to come:*

- ▶ **Practical uncertainty** in deep learning applications
- ▶ **Principled extensions** to deep learning tools
- ▶ **Hybrid** deep learning – Bayesian models

*and much, much, more.*

**Thank you for listening.**