

## 6. Variational Dropout

1) Gaussian Dropout:  $B = WA, W \sim \mathcal{N}(\hat{W}, 2 \odot \hat{W}^2)$   
 $W \sim \hat{W} \odot \mathcal{N}(1, 2)$

2) Variational Dropout:

ELBO:  $\mathcal{L}_V(\phi) = \mathbb{E}_{q(W|\phi)} \log p(y|x, W) - D_{KL}[q(W|\phi) \| p(W)] \rightarrow \max_{\phi}$

posterior:  $w_{ij} = \hat{w}_{ij} (1 + \sqrt{2} \epsilon_{ij}), \epsilon_{ij} \sim \mathcal{N}(0, 1)$

$q(w_{ij} | \phi_{ij}) = \mathcal{N}(w_{ij} | \hat{w}_{ij}, 2_{ij} \hat{w}_{ij}^2)$

prior:  $p(w_{ij}) \propto \frac{1}{|w_{ij}|}$  - log-uniform

KL-Divergence:  $-D_{KL}(q(w_{ij} | \hat{w}_{ij}, 2_{ij}) \| p(w_{ij})) =$

$= 0.5 \log 2_{ij} - \mathbb{E}_{\epsilon \sim \mathcal{N}(1, 2_{ij})} \log |\epsilon| + C$

$\uparrow$   
Doesn't depend on  $w_{ij}$ .  $\Rightarrow$  can fix 2 and optimize only  $\hat{w}_{ij}$

gaussian dropout is Bayesian technique

3) Local reparametrization to Bayesian Dropout:

$p(w_i) = \mathcal{N}(w_i | \mu_i, \sigma_i^2), b = w^T x = \sum_i (\mu_i + \epsilon_i \sigma_i) x_i = \sum_i \mu_i x_i + \epsilon \sqrt{\sum_i \sigma_i^2 x_i^2}$

$q(w_{ij} | \phi_{ij}) = \mathcal{N}(w_{ij} | \hat{w}_{ij}, 2_{ij} \hat{w}_{ij}^2)$

output =  $Wx$

output =  $\text{Relu}(\text{output})$

Output  $\sim \mathcal{N}(Wx, [W \odot W][x \odot x])$

Output =  $\text{Relu}(\text{Output})$

#### 4) Variance Reduction.

$$w_{ij} = \hat{w}_{ij} (1 + \sqrt{2_{ij}} \cdot \epsilon_{ij}) \Rightarrow \frac{\partial w_{ij}}{\partial \hat{w}_{ij}} = 1 + \sqrt{2_{ij}} \epsilon_{ij}$$

Very noisy for  $2 > 1$ .

Solution: Additive Noise Parametrization:

$$w_{ij} = \hat{w}_{ij} + \delta_{ij} \epsilon_{ij}, \quad \delta_{ij}^2 = 2_{ij} \hat{w}_{ij}^2$$

$$\frac{\partial w_{ij}}{\partial \hat{w}_{ij}} = 1 \leftarrow \text{no noise}$$

optimize the ELBO w.r.t.  $(\hat{w}, \delta)$

KL-div approximated analytically.

#### 5) ~~Warm-up~~ Pruning problem: ~~sol~~

Too many weights are pruned at the beginning

Use annealing of the KL coefficient with  $\tau$  weight.

Warmup solution.

One sigma per layer works just as well.

For large models - use pretrained model.

#### 6) Bayesian Jitterification of RNNs.

$$q(w_{ij}^x) = \mathcal{N}(w_{ij}^x | \hat{w}_{ij}^x, \delta_{ij}^{x^2}), \quad q(w_{ij}^h) = \mathcal{N}(w_{ij}^h | \hat{w}_{ij}^h, \delta_{ij}^{h^2})$$

- log-uniform prior

- same sample  $W$  for all timestamps

- cannot use local reparameteriz. trick

- up to 200x compression on ~~language modelling~~ sentiment analysis.

#### 7) Group Sparsity: Drop neurons instead of weights.

Cannot use one  $\delta$  per layer.

Log-Normal Dropout:  $p(\theta_i) = \text{Log } U_{[0,1]}(\theta_i) \quad q(\theta_i | \mu, \sigma_i) = \text{Log } \mathcal{N}_{[0,1]}(\theta_i | \mu, \sigma_i^2)$

(structured BP)