

# Bayesian Compression for Natural Language Processing

Nadezhda Chirkova<sup>\*1</sup> Ekaterina Lobacheva<sup>\*1</sup> Dmitry Vetrov<sup>1,2</sup>

<sup>1</sup>Higher School of Economics, Samsung-HSE Laboratory; <sup>2</sup>Samsung AI Center, Russia



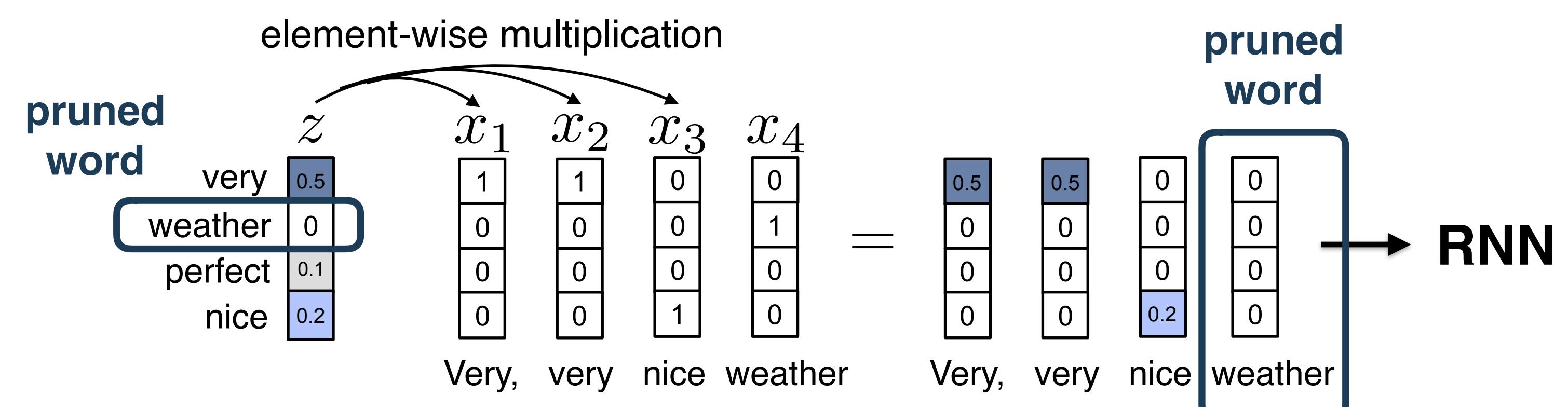
SAMSUNG  
Research

\* Equal contribution

Goal	Compress RNNs for NLP tasks: text classification and text generation
Method	<ul style="list-style-type: none"> <li>Apply Sparse Variational Dropout [1] taking into account RNN specifics</li> <li>Introduce group variables for elements of vocabulary</li> </ul>
Results	6–10K times weights compression and 2–70 times vocabulary compression depending on a task

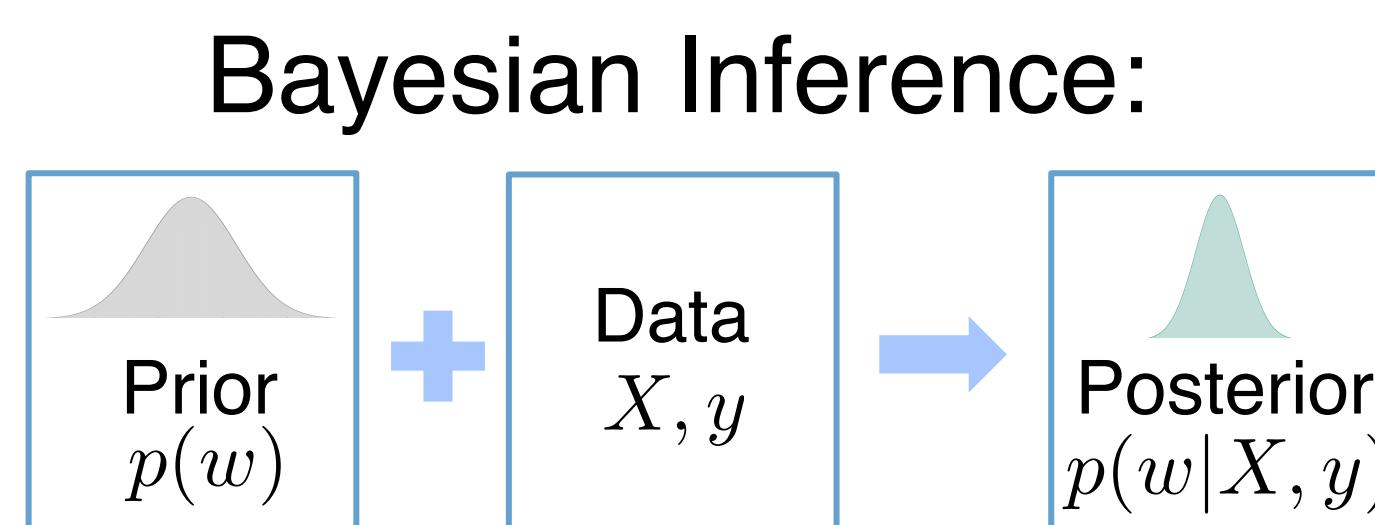
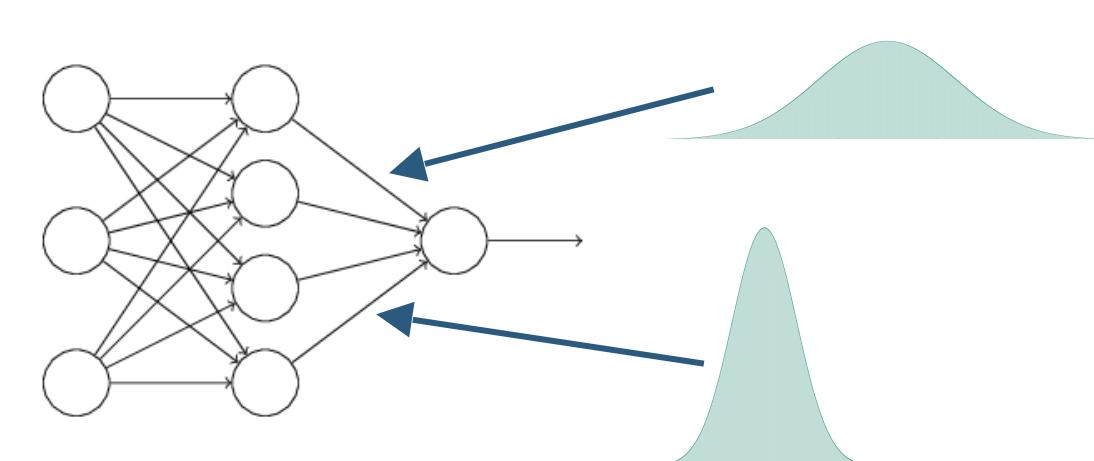
## Vocabulary sparsification

Input layer:  $x_t \rightarrow x_t \odot z$  where  $x_t$  is one-hot encoded token



Train and prune  $z$  in the same way as  $w$

## Bayesian neural networks



Posterior is intractable in NNs  $\rightarrow$  approximate it with  $q(w|\lambda)$ :

$$KL(q(w|\lambda)||p(w|X, y)) \rightarrow \min_{\lambda}$$

Equivalent to variational lower bound optimization:

$$\sum_{i=1}^N \underbrace{\mathbb{E}_{q(w|\lambda)} \log p(y^i|x^i, w)}_{\text{Data term}} - \underbrace{KL(q(w|\lambda)||p(w))}_{\text{Regularizer}} \rightarrow \max_{\lambda}$$

## Sparse Variational Dropout for RNNs

model Stochastic weights  $w$  + deterministic biases  $B$

Prior [2]:

$$p(w_{ij}) \propto 1/|w_{ij}|$$

Approximate posterior [2]:

$$q(w_{ij}|\theta_{ij}, \sigma_{ij}) = \mathcal{N}(w_{ij}|\theta_{ij}, \sigma_{ij}^2)$$

training Optimize ELBO to find  $\theta, \sigma, B$  [2]:

$$-\sum_{i=1}^N \int q(w|\theta, \sigma) \log p(y^i|x_0^i, \dots, x_T^i, w, B) dw + \sum_{w_{ij} \in \omega} KL(q(w_{ij}|\theta_{ij}, \sigma_{ij})||p(w_{ij})) \rightarrow \min_{\theta, \sigma, B}$$

data term  
sparsity inducing regularizer

• Integral estimation: 1 sample, reparametrization trick [2]:

$$w_{ij} = \theta_{ij} + \epsilon_{ij}\sigma_{ij}, \quad \epsilon_{ij} \sim \mathcal{N}(\epsilon_{ij}|0, 1)$$

• KL-divergence is approximated analytically [1]

After training:  $q(w_{ij}) = \delta(0)$  if  $\frac{\theta_{ij}}{\sigma_{ij}} <$  threshold

## RNN specifics

Recurrent layer:  $h_{t+1} = g(W^x x_{t+1} + W^h h_t + b)$

1 The same sample  $w$  for all  $t$  in one sequence

2 Local reparametrization trick (LRT) [2] is not applicable.

$$(Wx)_i = \sum_j \theta_{ij} x_j + \epsilon_i \sqrt{\sum_j (\sigma_{ij})^2 x_j^2}$$

normal r. v. constant 2d noise on  $W$   
↓  
1d noise on  $Wx$

In RNN:  $h_t$  is a random variable depending on  $W^h$   $\rightarrow$   $W^h h_t$  is not normal

Sampling the same noise on  $W^x$  for all t  $\neq$  Sampling the same noise on preactivation  $W^x x_t$  for all t

$$\sum_j \theta_{ij}^x x_{t,j} + \sum_j \epsilon_{ij}^x \sigma_{ij}^x x_{t,j} \neq \sum_j \theta_{ij}^x x_{t,j} + \epsilon_i^x \sqrt{\sum_j (\sigma_{ij}^x)^2 x_{t,j}^2}$$

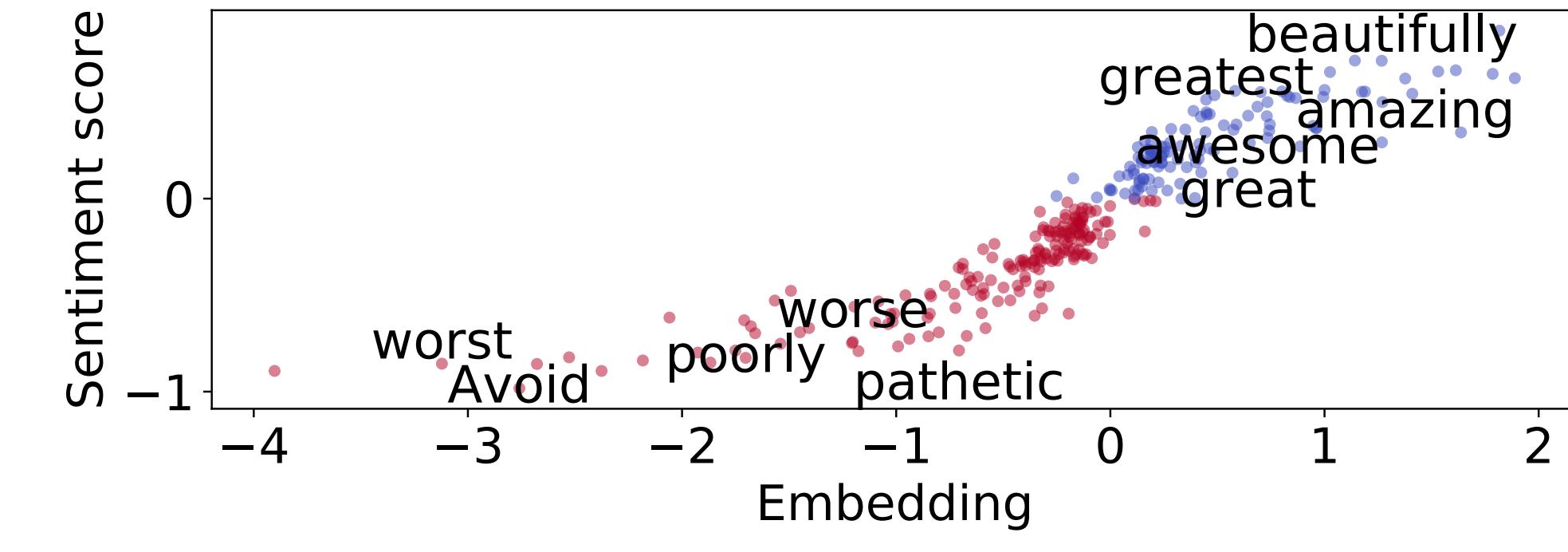
## Experiments: text classification

IMDb (25K texts, 2 classes), AGNews (120K texts, 4 classes)

Architecture: Embedding (300)  $\rightarrow$  LSTM (128 / 512)  $\rightarrow$  FC on  $h_T$

Task	Method	Accuracy %	Compression	Vocabulary
IMDb	Original	84.1	1x	20000
	SparseVD	<b>85.1</b>	1135x	4611
	SparseVD-Voc	83.6	<b>18792x</b>	<b>292</b>
AGNews	Original	<b>90.6</b>	1x	20000
	SparseVD	88.8	322x	5727
	SparseVD-Voc	89.2	<b>469x</b>	<b>2444</b>

IMDb: the only remaining embedding component for 292 remaining words is highly correlated with sentiment score



## Experiments: text generation

PTB Corpus: char-level (6M chars), word-level (0.9M words)

Architecture: LSTM (1000 / 256)  $\rightarrow$  FC

Task	Method	Valid	Test	Compression	Vocabulary
Char PTB (Bits-per-char)	Original	1.499	1.454	1x	50
	SparseVD	1.472	1.429	<b>7.9x</b>	50
	SparseVD-Voc	<b>1.458</b>	<b>1.417</b>	6.0x	<b>46</b>
Word PTB (Perplexity)	Original	135.6	129.3	1x	10000
	SparseVD	<b>115.0</b>	<b>109.2</b>	<b>22.1x</b>	9990
	SparseVD-Voc	126.0	120.2	19.3x	<b>3164</b>

### Paper



### Source Code



### Our follow-up on group sparsity for LSTMs



### References

- [1] Molchanov D., Ashukha A., Vetrov D. Variational dropout sparsifies deep neural networks. ICML 2017  
[2] Kingma D., Salimans T., Welling M. Variational dropout and the local reparameterization trick. NIPS 2015