# CS 6770 Natural Language Processing

Recurrent Neural Networks Language Models

Yangfeng Ji

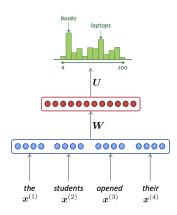
Information and Language Processing Lab Department of Computer Science University of Virginia

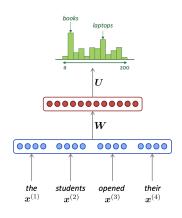


#### Overview

- 1. Neural Network Language Models
- 2. Recurrent Neural Networks
- 3. Computation Graph
- 4. RNN Language Modeling
- 5. Challenge of Training RNNs

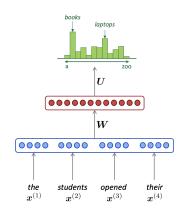
# Neural Network Language Models





Concatenated word embeddings

$$v = [v_{x_1}, v_{x_2}, v_{x_3}, v_{x_4}]$$
 (3)

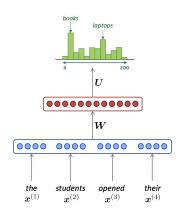


► Hidden layer:  $f(\cdot)$  could be any nonlinear activation function

$$h = f(Wv + b_1) \tag{2}$$

Concatenated word embeddings

$$v = [v_{x_1}, v_{x_2}, v_{x_3}, v_{x_4}]$$
 (3)



Output distribution

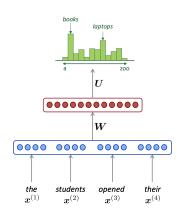
$$P(X_5 \mid X_{1:4}) = \operatorname{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|\mathcal{V}|}$$
(1)

► Hidden layer:  $f(\cdot)$  could be any nonlinear activation function

$$h = f(Wv + b_1) \tag{2}$$

Concatenated word embeddings

$$v = [v_{x_1}, v_{x_2}, v_{x_3}, v_{x_4}]$$
 (3)



Output distribution

$$P(X_5 \mid X_{1:4}) = \operatorname{softmax}(\boldsymbol{U}\boldsymbol{h} + \boldsymbol{b}_2) \in \mathbb{R}^{|\mathcal{V}|}$$
(1)

► Hidden layer:  $f(\cdot)$  could be any nonlinear activation function

$$h = f(Wv + b_1) \tag{2}$$

Concatenated word embeddings

$$v = [v_{x_1}, v_{x_2}, v_{x_3}, v_{x_4}]$$
 (3)

 $\blacktriangleright$  Word indices:  $x_1, x_2, x_3, x_4$ 

This is the very first neural neural language model [Bengio et al., 2001], which has a similar network architecture as the one discussed text classification.

# A Neural Probabilistic Language Model

The first paragraph of the paper *A Neural Probabilistic Language Model* [Bengio et al., 2001]

#### 1 Introduction

A fundamental problem that makes language modeling and other learning problems difficult is the *curse of dimensionality*. It is particularly obvious in the case when one wants to model the joint distribution between many discrete random variables (such as words in a sentence, or discrete attributes in a data-mining task). For example, if one wants to model the joint distribution of 10 consecutive words in a natural language with a vocabulary V of size 100,000, there are potentially  $100\,000^{10}-1=10^{50}-1$  free parameters.

<sup>&</sup>lt;sup>1</sup>For more precise description, please refer to [Shalev-Shwartz and Ben-David, 2014]

# A Neural Probabilistic Language Model

The first paragraph of the paper *A Neural Probabilistic Language Model* [Bengio et al., 2001]

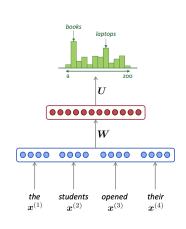
#### 1 Introduction

A fundamental problem that makes language modeling and other learning problems difficult is the *curse of dimensionality*. It is particularly obvious in the case when one wants to model the joint distribution between many discrete random variables (such as words in a sentence, or discrete attributes in a data-mining task). For example, if one wants to model the joint distribution of 10 consecutive words in a natural language with a vocabulary V of size 100,000, there are potentially  $100\,000^{10}-1=10^{50}-1$  free parameters.

#### Curse of dimensionality

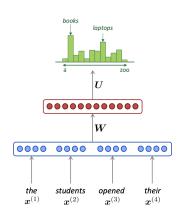
The sample complexity is an exponential function of the dimensionality of data<sup>1</sup>

<sup>&</sup>lt;sup>1</sup>For more precise description, please refer to [Shalev-Shwartz and Ben-David, 2014]



Improvement over *n*-gram language models

- ightharpoonup Less parameters (with large n's)
- No sparsity problem
- No smoothing is needed

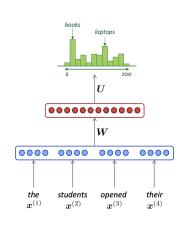


Improvement over *n*-gram language models

- ► Less parameters (with large *n*'s)
- No sparsity problem
- No smoothing is needed

#### Remaining issues

- Fixed window size similar to n-gram models
- No explicit modeling of word order beyond the window size
- ► Same word will be computed *k* times along the sliding window, where *k* is the window size



Improvement over *n*-gram language models

- ► Less parameters (with large *n*'s)
- No sparsity problem
- No smoothing is needed

#### Remaining issues

- Fixed window size similar to n-gram models
- No explicit modeling of word order beyond the window size
- ► Same word will be computed *k* times along the sliding window, where *k* is the window size

We need a new neural network architecture that can read words continuously along predictions

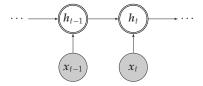
#### Recurrent Neural Networks

#### Recurrent Neural Networks (RNNs)

A simple RNN is defined by the following recursive function

$$h_t = f(x_t, h_{t-1}) \tag{4}$$

and depicted as



#### where

- ▶  $h_{t-1}$ : hidden state at time step t-1
- $\triangleright$   $x_t$ : input at time step t
- $\blacktriangleright$   $h_t$ : hidden state at time step t

# A Simple Transition Function

In the simplest case, the transition function f is defined with an element-wise Sigmoid function and a linear transformation of  $x_t$  and  $h_{t-1}$ 

$$h_t = f(x_t, h_{t-1}) = \sigma(W_h h_{t-1} + W_i x_t + b)$$
 (5)

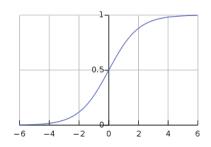
#### where

- $\triangleright$   $x_t$ : input word embedding
- ▶  $h_{t-1}$ : hidden statement from previous time step
- $\triangleright$  **W**<sub>h</sub>: parameter matrix for hidden states
- ► **W**<sub>i</sub>: parameter matrix for inputs
- ▶ *b*: bias term (also a parameter)

# Sigmoid Function

A Sigmoid function with one-dimensional input  $x \in (-\infty, \infty)$ 

$$\sigma(x) = \frac{1}{1 - e^{-x}}$$



The potential numeric issue caused by the Sigmoid function

- $\sigma(x) \to 1 \text{ with } x \gg 6$
- $ightharpoonup \sigma(x) o 0, x \ll -6$

The output of the Sigmoid function will approximate a constant, when the input value is beyond certain ranges

# **Unfolding RNNs**

We can unfold this recursive definition of a RNN

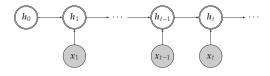
$$h_t = f(x_t, h_{t-1}) \tag{6}$$

# **Unfolding RNNs**

We can unfold this recursive definition of a RNN

$$h_t = f(x_t, h_{t-1}) \tag{6}$$

as



$$h_{t} = f(x_{t}, f(x_{t-1}, h_{t-2}))$$

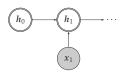
$$= f(x_{t}, f(x_{t-1}, f(x_{t-2}, h_{t-3})))$$

$$= \cdots$$

$$= f(x_{t}, f(x_{t-1}, f(x_{t-2}, \cdots f(x_{1}, h_{0}) \cdots)))$$
 (7)

#### **Base Condition**

Base condition defines the starting point of the recursive computation



$$h_t = f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(x_1, h_0) \dots)))$$
 (8)

- $\blacktriangleright$   $h_0$ : zero vector or parameter
- $\triangleright$   $x_1$ : input at time t=1

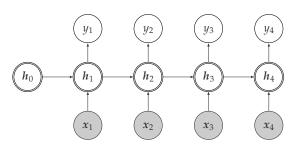
11

#### RNN for Sequential Prediction

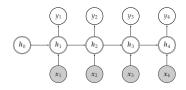
In general, RNNs can be used for any sequential modeling tasks

$$h_t = f(x_t, h_{t-1}) (9)$$

$$P(Y_t; h_t) = \operatorname{softmax}(W_o h_t + b_o)$$
 (10)



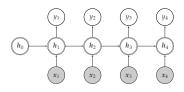
# Sequential Modeling as Classification



▶ Prediction at each time step *t* 

$$\hat{y}_t = \operatorname*{argmax}_{y} P(Y_t = y; h_t)$$
 (11)

# Sequential Modeling as Classification



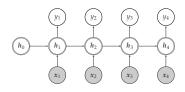
Prediction at each time step t

$$\hat{y}_t = \underset{y}{\operatorname{argmax}} P(Y_t = y; h_t)$$
 (11)

Loss at single time step t

$$L_t(y_t, \hat{y}_t) = -\log P(y_t; h_t)$$
(12)

# Sequential Modeling as Classification



▶ Prediction at each time step *t* 

$$\hat{y}_t = \underset{y}{\operatorname{argmax}} P(Y_t = y; \boldsymbol{h}_t)$$
 (11)

Loss at single time step t

$$L_t(y_t, \hat{y}_t) = -\log P(y_t; h_t)$$
(12)

► The total loss

$$\ell = \sum_{t=1}^{T} L_t(y_t, \hat{y}_t)$$
 (13)

# Computation Graph

#### **Forward Operations**

For simplicity, consider the example of a two-layer neural network

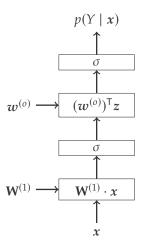
$$P(Y = +1 \mid \mathbf{x}) = \sigma\left((\mathbf{w}^{(o)})^{\mathsf{T}}\sigma(\mathbf{W}^{(1)}\mathbf{x})\right)$$
(14)

A neural network is a composition of some basic functions and operations. For example

- the Sigmoid function  $\sigma(\cdot)$
- matrix transpose  $(w^{(o)})^T$
- ightharpoonup matrix-vector multiplication  $\mathbf{W}^{(1)}x$

# Forward Graph

The computation graph of the two-layer neural network<sup>2</sup>



<sup>&</sup>lt;sup>2</sup>For simplicity, the *transpose* operation blocks are ignored from the graph

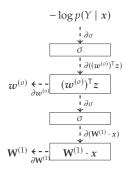
#### **Backward Operations**

Similarly, the gradient of neural network parameters are computed with a series of backward operations associated with the derivative of some basic function. For example

$$ightharpoonup \frac{\mathrm{d}bx}{\mathrm{d}x} = b$$

# **Backward Graph**

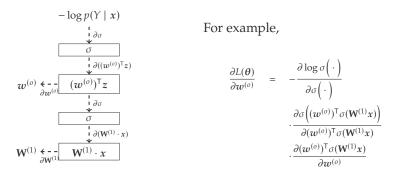
- ► The gradient of each operations can be computed individually based on the input and output
- With the chain rule, gradient of the loss function with respect to any parameter is a sequence of multiplications of individual gradients



x

# **Backward Graph**

- ► The gradient of each operations can be computed individually based on the input and output
- With the chain rule, gradient of the loss function with respect to any parameter is a sequence of multiplications of individual gradients



x

#### Example: MiniTorch

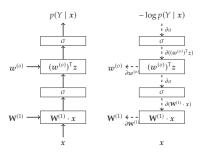
#### A screenshot from the MiniTorch library<sup>3</sup>

```
class Add(Function):
    @staticmethod
    def forward(ctx, t1, t2):
        return add zip(t1, t2)
    @staticmethod
    def backward(ctx, grad output):
        return grad output, grad output
class Mul(Function):
    @staticmethod
    def forward(ctx, a, b):
        ctx.save for backward(a, b)
        return mul zip(a, b)
    @staticmethod
    def backward(ctx, grad output):
        a, b = ctx.saved values
        return mul zip(b, grad output), mul zip(a, grad output)
class Sigmoid(Function):
    @staticmethod
    def forward(ctx, a):
        out = sigmoid map(a)
        ctx.save for backward(out)
        return out
    @staticmethod
    def backward(ctx, grad output):
        sigma = ctx.saved values
        return sigma * (-sigma + 1.0) * grad output
```

<sup>&</sup>lt;sup>3</sup>A teaching library of PyTorch developed by Sasha Rush and Ge Gao

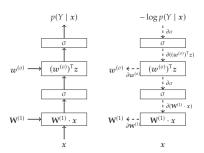
# **Computation Graph**

Perform the forward/backward step with a graph of basic operations (e.g., PyTorch, Tensorflow)



#### **Computation Graph**

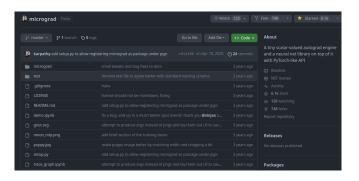
Perform the forward/backward step with a graph of basic operations (e.g., PyTorch, Tensorflow)



- Modular implementation: implement each module with its forward/backward operations together
- Automatic differentiation: automatically run with the backward step

#### Example: Micrograd

MICROGRAD is an extremely simple example to illustrate the idea of using autograd to implementing the backpropagation algorithm.



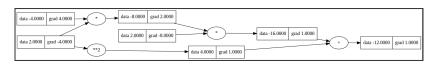
#### Demo

# **Example: Gradient Computation**

**Function** 

$$f(a,b) = 2ab + b^2$$

The computation graph on forward and backward operations



RNN Language Modeling

## Language Models

A language model defines the probability of  $x_t$  given  $x = (x_1, x_2, ..., x_{t-1})$  as

$$P(x_t \mid x_1, \dots, x_{t-1}) \tag{15}$$

and the joint probability as

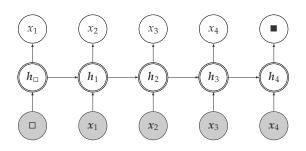
$$P(\mathbf{x}_{1:T}) = P(x_1) \cdot P(x_2 \mid x_1)$$

$$\cdot \cdot \cdot \cdot \cdot$$

$$\cdot P(x_T \mid x_1, x_2, \dots, x_{T-1})$$

#### Language Modeling with RNNs

Using RNNs for language modeling

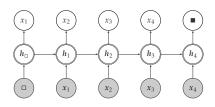


with two special tokens

$$\{\Box, x_1, \ldots, x_T, \blacksquare\}$$

#### **RNN Language Models**

For a given sentence  $\{x_1, \ldots, x_t\}$ , the input at time t is word embedding  $x_t$ 



The probability distribution of next word  $X_t$ 

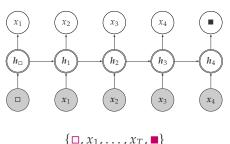
$$P(X_t = x \mid \mathbf{x}_{1:t-1}) = \frac{\exp(\mathbf{w}_{o,x}^{\mathsf{T}} \mathbf{h}_{t-1})}{\sum_{x' \in \mathcal{V}} \exp(\mathbf{w}_{o,x'}^{\mathsf{T}} \mathbf{h}_{t-1})}$$
(16)

where

- $w_{o,x}$  is the output weight vector (parameter) associated with word x
- $ightharpoonup \mathscr{V}$  is the word vocabulary

## **Special Cases**

Similar to statistical language modeling, there are also two special cases that we need to consider



The corresponding prediction functions are defined as

ightharpoonup At time t = 1

$$P(X_1 = x) \propto \exp(\boldsymbol{w}_{o,x}^{\mathsf{T}} \boldsymbol{h}_{\square}) \tag{17}$$

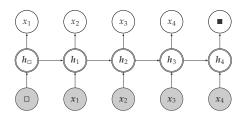
At time t = T

$$P(X_T = \blacksquare \mid x_{1:T-1}) \propto \exp(w_{o,x}^\mathsf{T} h_{T-1})$$
(18)

# Challenge of Training RNNs

#### Objective

The training objective for each timestep is to predict the next token in the text

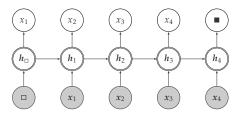


- ► Prediction at step t,  $P(X_t = x \mid x_{1:t-1}) = \frac{\exp(w_{o,x}^\mathsf{T} h_{t-1})}{\sum_{x' \in \mathscr{V}} \exp(w_{o,x'}^\mathsf{T} h_{t-1})}$
- Loss at step t,  $L_t = -\log P(X_t = x \mid x_{1:t-1})$

#### Gradients

Let  $\theta$  denote all model parameters

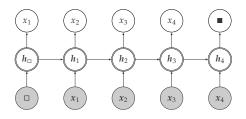
$$\frac{\partial \ell}{\partial \boldsymbol{\theta}} = \sum_{t=1}^{T} \frac{\partial L_t}{\partial \boldsymbol{\theta}} \tag{19}$$



Backpropagation Through Time [Rumelhart et al., 1985, BPTT]

#### **Model Parameters**

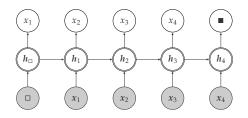
Before computing the gradient of each  $L_t$  with respect to model parameters, let us count how many parameters that we need consider



• Output parameter matrix  $W_0 = (w_{0,1}, \dots, w_{0,V})$ 

#### **Model Parameters**

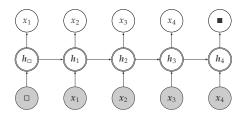
Before computing the gradient of each  $L_t$  with respect to model parameters, let us count how many parameters that we need consider



- Output parameter matrix  $W_o = (w_{o,1}, \dots, w_{o,V})$
- ► Input word embedding matrix  $X = (x_1, ..., x_V)$

#### **Model Parameters**

Before computing the gradient of each  $L_t$  with respect to model parameters, let us count how many parameters that we need consider



- Output parameter matrix  $W_o = (w_{o,1}, \dots, w_{o,V})$
- ► Input word embedding matrix  $X = (x_1, ..., x_V)$
- Neural network parameters  $W_h$ ,  $W_i$ , b

#### **Backpropagation Through Time**

Take time step t as an example, we can take a look the gradient computation of some specific parameters

▶ Output model parameter  $\frac{\partial L_t}{\partial w_{o,\cdot}}$ 

#### **Backpropagation Through Time**

Take time step t as an example, we can take a look the gradient computation of some specific parameters

- Output model parameter  $\frac{\partial L_t}{\partial w_{o,\cdot}}$
- Neural network parameters, for example W<sub>h</sub>

$$\frac{\partial L_t}{\partial W_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial h_t} \cdot \left( \prod_{j=i}^{t-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_i}{\partial W_h} \right\} \tag{20}$$

Similar patterns for the other two neural network parameters  $W_i$  and b

## **Backpropagation Through Time**

Take time step t as an example, we can take a look the gradient computation of some specific parameters

- Output model parameter  $\frac{\partial L_t}{\partial w_{o,\cdot}}$
- ▶ Neural network parameters, for example  $W_h$

$$\frac{\partial L_t}{\partial W_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial h_t} \cdot \left( \prod_{j=i}^{t-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_i}{\partial W_h} \right\} \tag{20}$$

Similar patterns for the other two neural network parameters  $W_i$  and b

- ▶ Word embedding  $\frac{\partial L_t}{\partial x_{t'}}$ 
  - ► E.g., word embedding  $x_{t'}$  is the input of  $h_t$  if  $t' \le t$ , so ...

## Challenges

For each timestep, we need to compute the gradient using the chain rule:

$$\frac{\partial L_t}{\partial W_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial h_t} \cdot \left( \prod_{j=i}^{t-1} \frac{\partial h_{j+1}}{\partial h_j} \right) \cdot \frac{\partial h_i}{\partial W_h} \right\}$$
(21)

The chain rule of gradient will cause two potential problems in training RNNs

- ▶ vanishing gradients:  $\frac{\partial L_t}{\partial \theta} \rightarrow 0$
- exploding gradients:  $\frac{\partial L_t}{\partial \theta} \ge M$

[Pascanu et al., 2013]

## **Exploding Gradients**

Solution: **norm clipping** [Pascanu et al., 2013].

Consider the gradient  $g = \frac{\partial \ell}{\partial \theta}$ ,

$$\hat{g} \leftarrow \tau \cdot \frac{g}{\|g\|} \tag{22}$$

when  $\|g\| > \tau$ .

- Usually,  $\tau = 3$  or 5 in practice.
- Smaller gradient will cause slower learning progress

## Vanishing Gradients

#### Solution:

- ▶ initialize parameters carefully
- replace hidden state transition function  $\sigma(\cdot)$  with other options

$$f(x_t, h_{t-1}) = \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_i x_t + b)$$
 (23)

- ► LSTM [Hochreiter and Schmidhuber, 1997]
- ► GRU [Cho et al., 2014]

## **Long Short-Term Memory**

## From the first page of the original paper proposing LSTM [Hochreiter and Schmidhuber, 1997]

The problem. With conventional "Back-Propagation Through Time" (BPTT, e.g., Williams and Zipser 1992, Werbos 1988) or "Real-Time Recurrent Learning" (RTRL, e.g., Robinson and Fallside 1987), error signals "flowing backwards in time" tend to either (1) blow up or (2) vanish: the temporal evolution of the backpropagated error exponentially depends on the size of the weights (Hochreiter 1991). Case (1) may lead to oscillating weights, while in case (2) learning to bridge long time lags takes a prohibitive amount of time, or does not work at all (see section 3).

The remedy. This paper presents "Long Short-Term Memory" (LSTM), a novel recurrent network architecture in conjunction with an appropriate gradient-based learning algorithm. LSTM is designed to overcome these error back-flow problems. It can learn to bridge time intervals in excess of 1000 steps even in case of noisy, incompressible input sequences, without loss of short time lag capabilities. This is achieved by an efficient, gradient-based algorithm for an architecture

## Long Short-Term Memory

Rather than directly taking input and hidden state as simple transition function, LSTM relies on three cates to control *how much* information it should take from input and hidden state before combining them together

$$\begin{array}{lll} i_t & = & \sigma(\mathbf{W}_{xi}x_t + \mathbf{W}_{hi}h_{t-1} + \mathbf{W}_{ci}c_{t-1} + b_i) \\ f_t & = & \sigma(\mathbf{W}_{xf}x_t + \mathbf{W}_{hf}h_{t-1} + \mathbf{W}_{cf}c_{t-1} + b_f) \\ c_t & = & f_t \circ c_{t-1} + i_t \circ \tanh(\mathbf{W}_{xc}x_t + \mathbf{W}_{hc}h_{t-1} + b_c) \\ o_t & = & \sigma(\mathbf{W}_{xo}x_t + \mathbf{W}_{ho}h_{t-1} + \mathbf{W}_{co}c_t + b_o) \\ h_t & = & o_t \circ \tanh(c_t) \end{array}$$

where  $\circ$  is the element-wise multiplication,  $\{W_\cdot\}$  and  $\{b_\cdot\}$  are parameters. [Graves, 2013]

#### Reference



Bengio, Y., Ducharme, R., and Vincent, P. (2001).

A neural probabilistic language model. In  $\it NIPS$ .



Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014).

On the properties of neural machine translation: Encoder-decoder approaches. arXiv preprint arXiv:1409.1259.



Graves, A. (2013).

Generating sequences with recurrent neural networks. arXiv preprint arXiv:1308.0850.



Hochreiter, S. and Schmidhuber, J. (1997).

Long short-term memory.

Neural computation, 9(8):1735-1780.



Pascanu, R., Mikolov, T., and Bengio, Y. (2013).

On the difficulty of training recurrent neural networks.

In International Conference on Machine Learning, pages 1310–1318.



Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985).

Learning internal representations by error propagation.

Technical report, California Univ San Diego La Jolla Inst for Cognitive Science.



Shalev-Shwartz, S. and Ben-David, S. (2014).

Understanding machine learning: From theory to algorithms.

Cambridge university press.