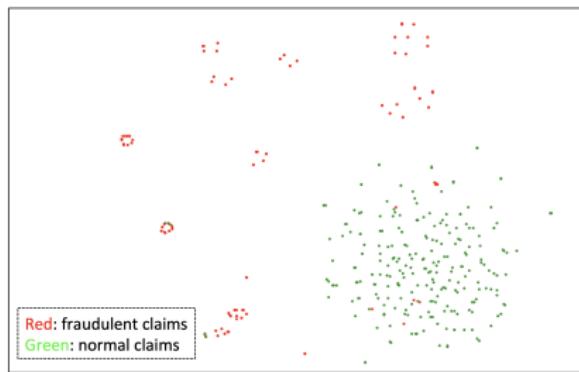
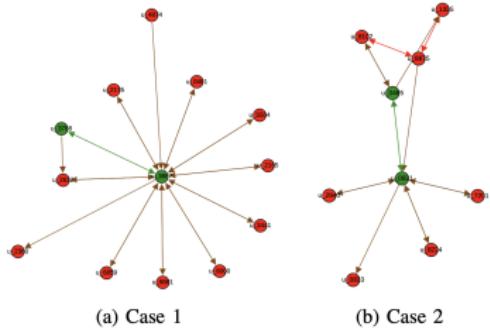


Intro to Graph Representation Learning

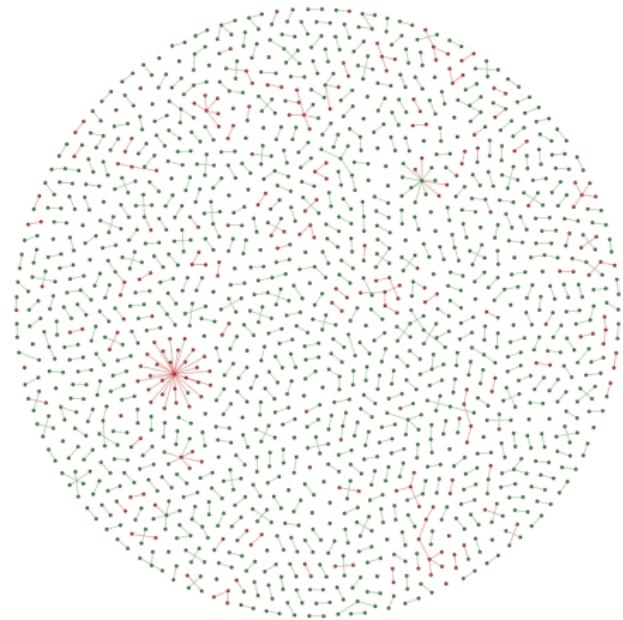
Dmitry Eremeev

Moscow, 2021

- Fraudulent claims in insurance.
- Graph of client calculations and transactions.



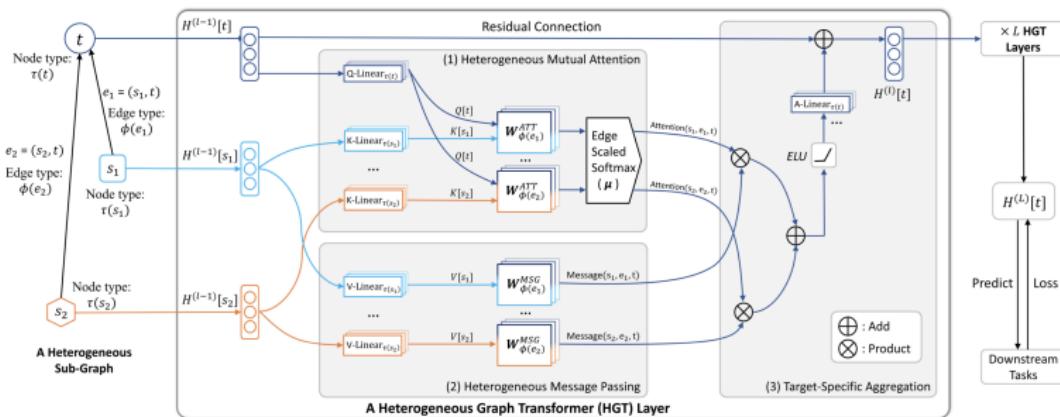
[source](#)



[source](#)

Heterogeneous Graph Transformer

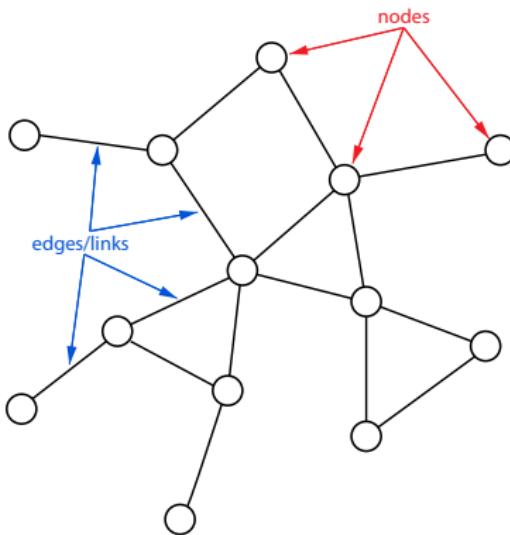
arXiv:2003.01332 [cs.LG]

github.com/acbull/pyHGT

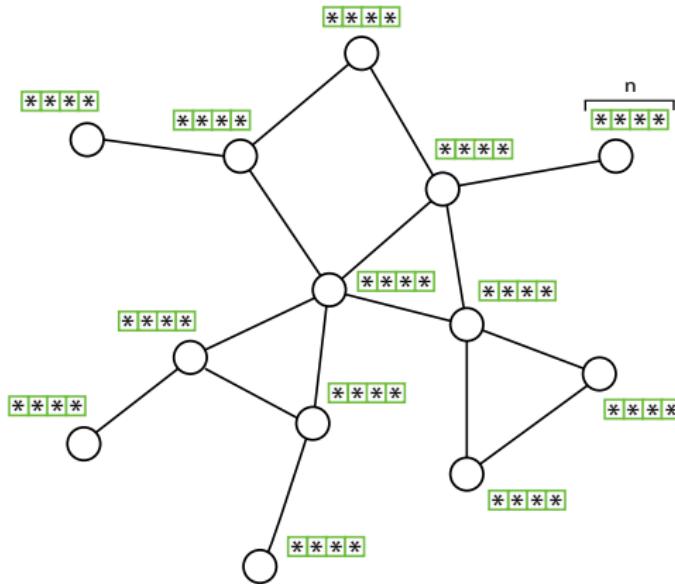
Intro

Graph $G = (V, E)$ is a pair of 2 sets:

- V - vertices [nodes],
- $E \in V \times V$ - edges [links] (paired vertices).

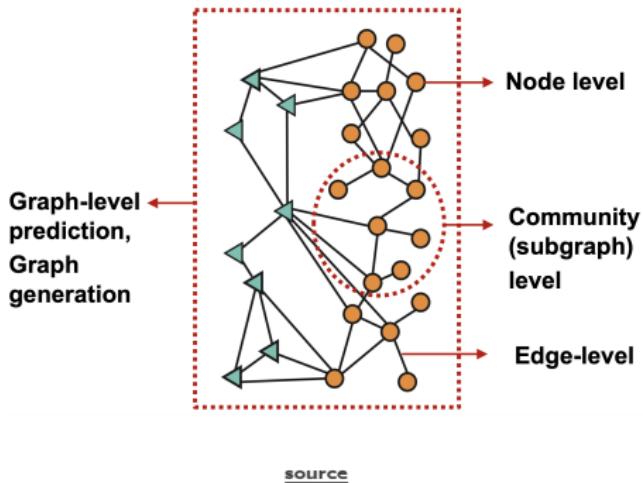


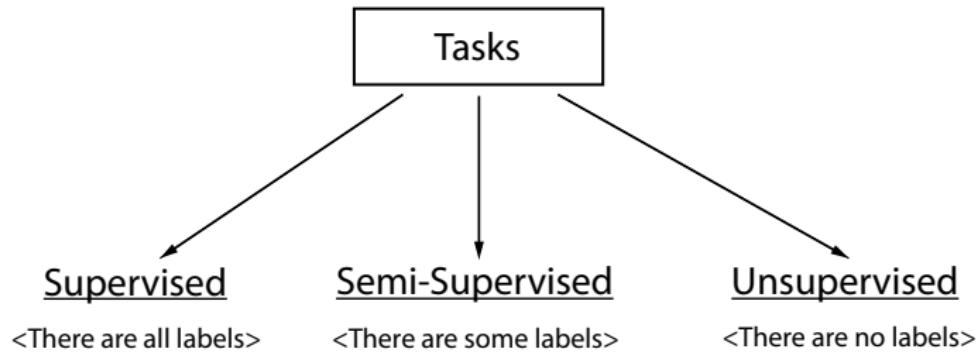
Each node could be described by feature vectors



Types of tasks:

- Node Classification
- Link prediction
- Community detection
- Learning over the whole graph



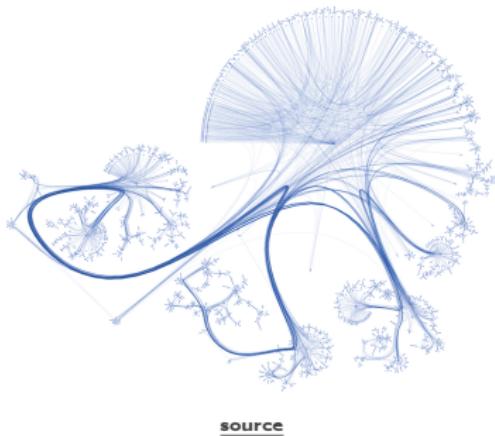


Networkx library

- „Hand-Crafted“ features.
- Node Degree, Node Centrality, Node Clustering Coefficient, etc.
- Check the Stanford CS224W course.

Cora dataset

- 2708 scientific publications, classified into 7 classes.
- 5429 links - citations.
- Features: presence of the words from dictionary.



Open Graph Benchmark


[Get Started](#) [Updates](#) [KDD Cup](#) [Datasets](#) ▾ [Leaderboards](#) ▾ [Paper](#) [Team](#) [Github](#)

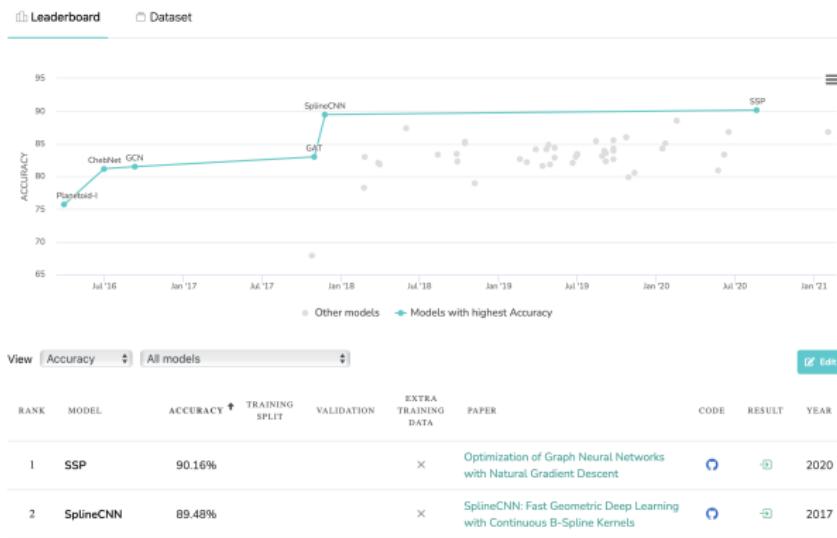
The classification accuracy on the test and validation sets. The higher, the better.

Package: >=1.1.1

Rank	Method	Test Accuracy	Validation Accuracy	Contact	References	#Params	Hardware	Date
1	MLP + C&S	0.8418 ± 0.0007	0.9147 ± 0.0009	Horace He (Cornell)	Paper , Code	96,247	GeForce RTX 2080 (11GB GPU)	Oct 27, 2020
2	Linear + C&S	0.8301 ± 0.0001	0.9134 ± 0.0001	Horace He (Cornell)	Paper , Code	10,763	GeForce RTX 2080 (11GB GPU)	Oct 27, 2020
3	UniMP	0.8256 ± 0.0031	0.9308 ± 0.0017	Yunsheng Shi (PGL team)	Paper , Code	1,475,605	Tesla V100 (32GB)	Sep 8, 2020
4	Plain Linear + C&S	0.8254 ± 0.0003	0.9103 ± 0.0001	Horace He (Cornell)	Paper , Code	4,747	GeForce RTX 2080 (11GB GPU)	Oct 27, 2020
5	DeeperGCN+FLAG	0.8193 ± 0.0031	0.9221 ± 0.0037	Kezhi Kong	Paper , Code	253,743	NVIDIA Tesla V100 (32GB GPU)	Oct 20, 2020
6	GAT+FLAG	0.8176 ± 0.0045	0.9251 ± 0.0006	Kezhi Kong	Paper , Code	751,574	GeForce RTX 2080 Ti (11GB GPU)	Oct 20, 2020

[paperswithcode](#)

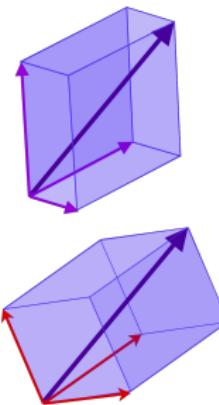
Node Classification on Cora



- Vectors: $\mathbf{v} \in V$, V - vector space.
- Basis: $B = \{\mathbf{b}_i\}$ - set of vectors that spans the whole vector space V :

$$\mathbf{v} = a_1 \mathbf{b}_1 + a_2 \mathbf{b}_2 + \cdots + a_n \mathbf{b}_n$$

$$\mathbf{v} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \in \mathbb{R}^n$$



- Matrix \mathbf{A} with m rows and n columns:

$$\mathbf{A} = \begin{matrix} & \begin{matrix} 1 & 2 & \cdots & n \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ m \end{matrix} & \left[\begin{matrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \cdots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{matrix} \right] \end{matrix} \in \mathbb{R}^{m \times n}$$

- Matrix as a set of vectors:

$$\mathbf{A} = \left[\begin{matrix} \text{---} & \mathbf{a}_1 & \text{---} \\ \text{---} & \mathbf{a}_2 & \text{---} \\ \vdots & & \vdots \\ \text{---} & \mathbf{a}_m & \text{---} \end{matrix} \right]$$

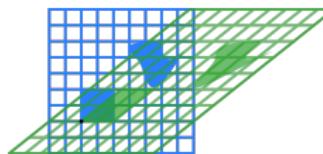
- Matrices define linear transformations:

$$\mathbf{A} \in \mathbb{R}^{m \times n} \iff \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\mathbf{v} \mapsto \mathbf{v}' = \mathbf{A}\mathbf{v}.$$

- Examples: $\mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\mathbf{A} = \begin{bmatrix} 1 & 1.25 \\ 0 & 1 \end{bmatrix}$$



$$\mathbf{A} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$$



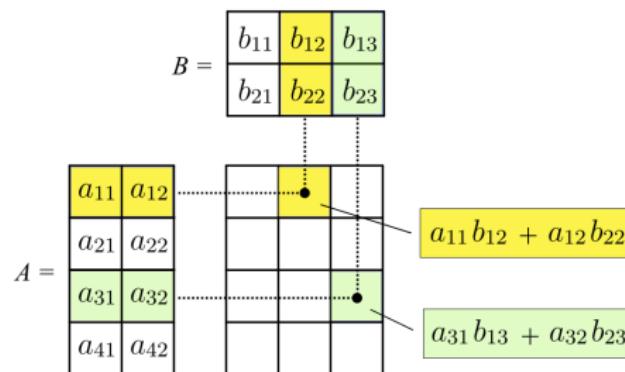
- Sum of matrices:

$$\mathbf{A} + \mathbf{B} := \begin{matrix} & 1 & \dots & n \\ \begin{matrix} 1 \\ 2 \\ \vdots \\ m \end{matrix} & \left[\begin{matrix} a_{11} + b_{11} & \dots & a_{1n} + b_{1n} \\ a_{21} + b_{21} & \dots & a_{2n} + b_{2n} \\ \vdots & \dots & \vdots \\ a_{m1} + b_{m1} & \dots & a_{mn} + b_{mn} \end{matrix} \right] \end{matrix}$$

- Matrix multiplication:

$$\underbrace{\mathbf{A}}_{m \times p} \cdot \underbrace{\mathbf{B}}_{p \times n} = \underbrace{\mathbf{C}}_{m \times n}$$

- Elements of \mathbf{C} are defined by **scalar products** of \mathbf{A} rows and \mathbf{B} columns:



27.3. A tensor of type (p, q) on V is an element of the space

$$T_p^q(V) = \underbrace{V^* \otimes \cdots \otimes V^*}_{p \text{ factors}} \otimes \underbrace{V \otimes \cdots \otimes V}_{q \text{ factors}}$$

isomorphic to the space of linear functions on $V \times \cdots \times V \times V^* \times \cdots \times V^*$ (with p factors V and q factors V^*). The number p is called the *covariant valency* of the tensor, q its *contravariant valency* and $p + q$ its *total valency*. The *vectors* are tensors of type $(0, 1)$ and *covectors* are tensors of type $(1, 0)$.

Let a basis e_1, \dots, e_n be selected in V and let e_1^*, \dots, e_n^* be the dual basis of V^* . Each tensor T of type (p, q) is of the form

$$T = \sum T_{i_1 \dots i_p}^{j_1 \dots j_q} e_{i_1}^* \otimes \cdots \otimes e_{i_p}^* \otimes e_{j_1} \otimes \cdots \otimes e_{j_q}; \quad (1)$$

the numbers $T_{i_1 \dots i_p}^{j_1 \dots j_q}$ are called the *coordinates of the tensor T* in the basis e_1, \dots, e_n .

Let us establish how coordinates of a tensor change under the passage to another basis. Let $\varepsilon_j = Ae_j = \sum a_{ij} e_i$ and $\varepsilon_j^* = \sum b_{ij} e_i^*$. It is easy to see that $B = (A^T)^{-1}$, cf. 5.3.

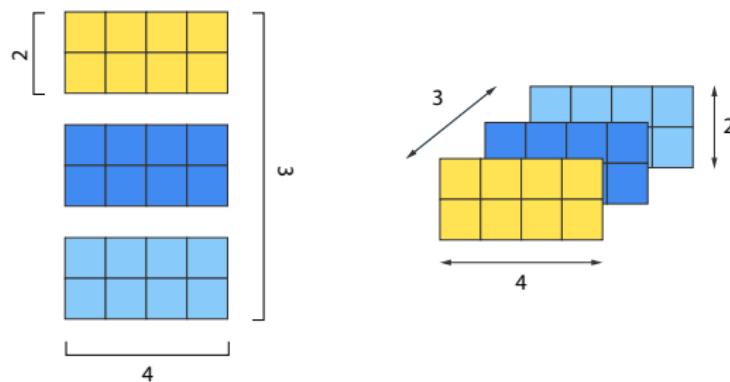
Introduce notations: $a_j^i = a_{ij}$ and $b_i^j = b_{ij}$ and denote the tensor (1) by $\sum T_\alpha^\beta e_\alpha^* \otimes e_\beta$ for brevity. Then

$$\sum T_\alpha^\beta e_\alpha^* \otimes e_\beta = \sum S_\mu^\nu \varepsilon_\mu^* \otimes \varepsilon_\nu = \sum S_\mu^\nu b_\mu^\alpha a_\nu^\beta e_\alpha^* \otimes e_\beta,$$

i.e.,

$$T_{i_1 \dots i_p}^{j_1 \dots j_q} = b_{i_1}^{l_1} \dots b_{i_p}^{l_p} a_{k_1}^{j_1} \dots a_{k_q}^{j_q} S_{l_1 \dots l_p}^{k_1 \dots k_q} \quad (2)$$

- Tensor of shape $[3, 2, 4]$ (a set of matrices $\mathbb{R}^{2 \times 4}$):





- PyTorch: pytorch.org
- PyTorch Geometric: github.com/rusty1s/pytorch_geometric

Tensor multiplication in PyTorch:

```
In 1  1 import torch

In 2  1 a = torch.tensor([[1, 2],
2   [3, 4]])
3 print(a.shape)
4 a

Out 2      torch.Size([2, 2])
            tensor([[1, 2],
               [3, 4]])

In 3  1 b = torch.tensor([[2, 0],
2   [0, 2]])
3 print(b.shape)
4 b

Out 3      torch.Size([2, 2])
            tensor([[2, 0],
               [0, 2]])

In 4  1 torch.matmul(a, b)

Out 4      tensor([[2, 4],
               [6, 8]])
```

PyTorch supports broadcasting:

```
In 5  1 # tensor of shape (2, 2, 2)
2 a = torch.randint(0, 10, (2,2,2))
3 print(a.shape)
4 a
```

```
Out 5      torch.Size([2, 2, 2])
            tensor([[[0, 3],
                         [7, 9]],
                         [[4, 6],
                         [5, 0]]])
```

```
In 10 1 b = torch.tensor([[2, 0, 1],
                           [0, 2, 1]])
2
3
4 # broadcast: [2, 2, 2] x [2, 3] --> [2, 2, 3]
5 c = torch.matmul(a, b)
6 print(c.shape)
7 c
```

```
Out 10      torch.Size([2, 2, 3])
            tensor([[[ 8, 16, 12],
                         [10, 16, 13]],
                         [[18,  4, 11],
                         [ 8, 14, 11]]])
```

Computing gradients:

```
In 1 1 import torch

In 2 1 x = torch.tensor(10.0, requires_grad=False)
2 y = torch.tensor(5.0, requires_grad=True)

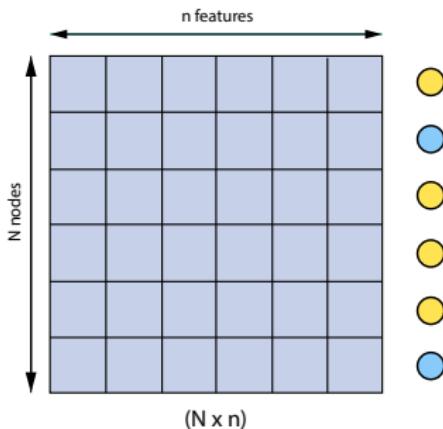
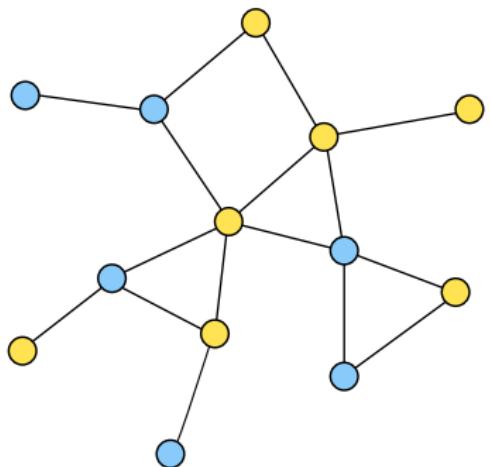
In 3 1 print(x)
2 print(y)

tensor(10.)
tensor(5., requires_grad=True)

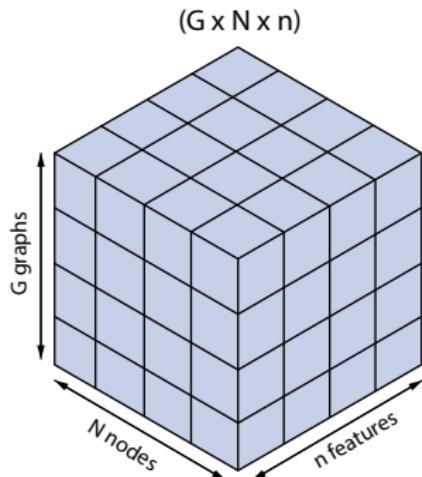
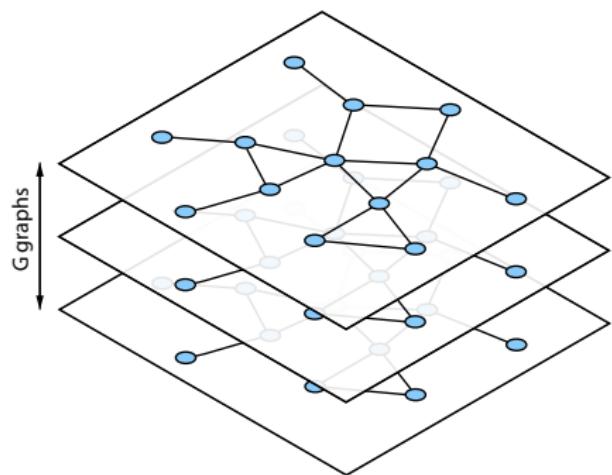
In 4 1 # compute the gradient
2 # z(x, y) = x*y^2
3 z = x*y*y
4 z.backward() # Autograd calculates and stores the gradient of z w.r.t to tensors
5 y.grad # extract the value of gradient

Out 4 tensor(100.)
```

Matrix representation

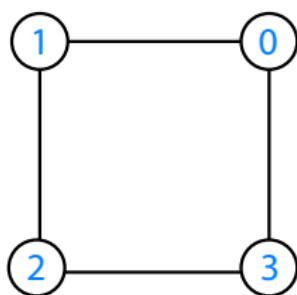


Tensor representation



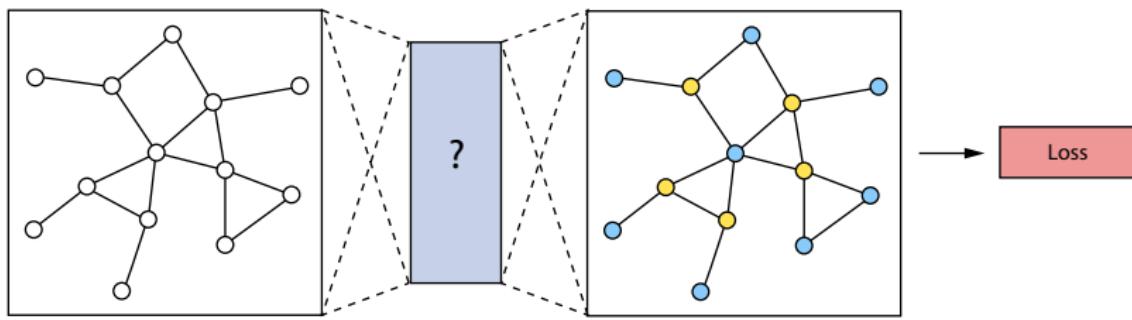
Learning over the whole graphs

- Contains information whether pairs of nodes are adjacent or not.
- Handy for determining neighbours for each node and creating masks.

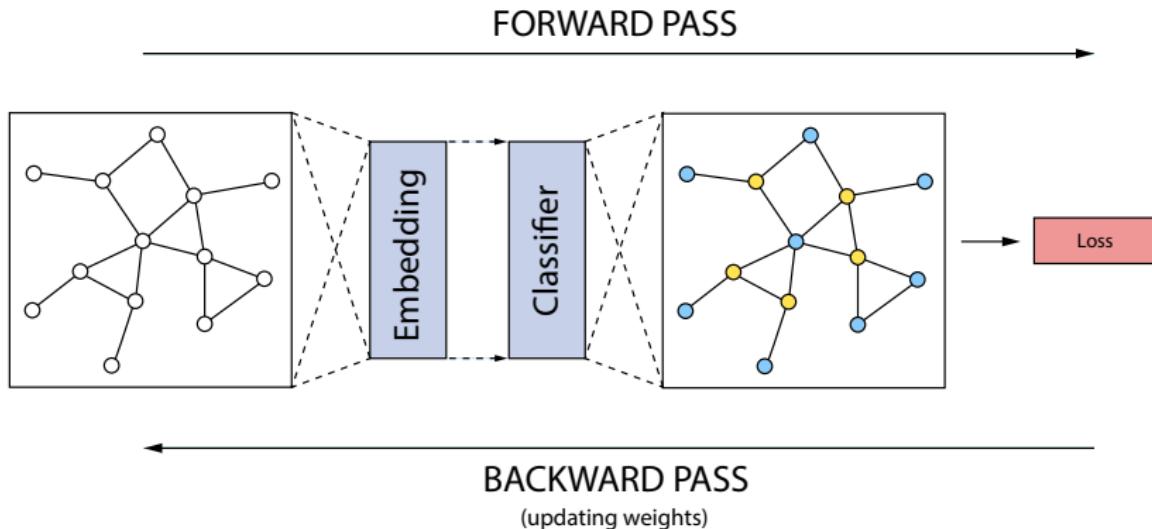


	0	1	2	3
0	0	1	0	1
1	1	0	1	0
2	0	1	0	1
3	1	0	1	0

FORWARD PASS

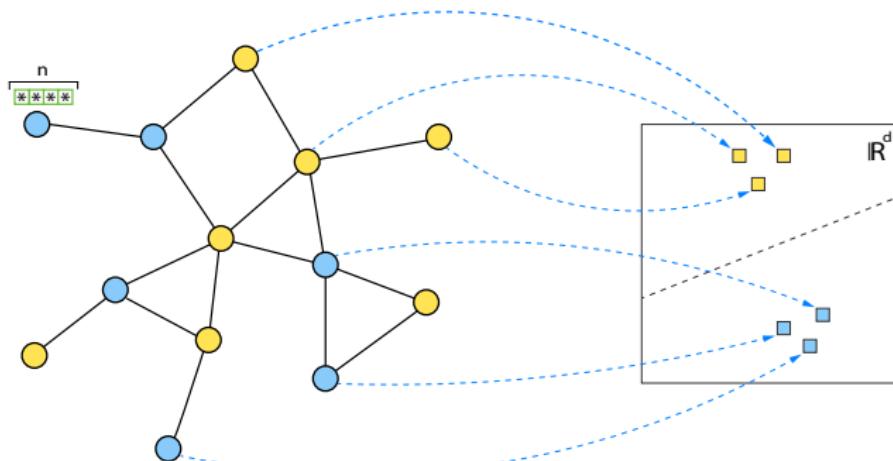


BACKWARD PASS (updating weights)



- Embeddings - learnable vector representations of nodes.
- Embeddings extract structural features, taking into account information in some neighbourhood.
- **Similarity** of nodes u and v is defined by **scalar product** of their embeddings $\mathbf{z}_u, \mathbf{z}_v$:

$$\text{Similarity}(u, v) \sim \langle \mathbf{z}_u; \mathbf{z}_v \rangle$$



Model Preliminaries

- Assume each point Y_i in data \mathcal{D} comes from exponential probability distribution p defined by parameters θ , ϕ and some functions a, c :

$$Y_i \sim p(y_i|\theta_i, \phi) = c(y_i, \phi) \exp\left(\frac{y_i\theta_i - a(\theta_i)}{\phi}\right).$$

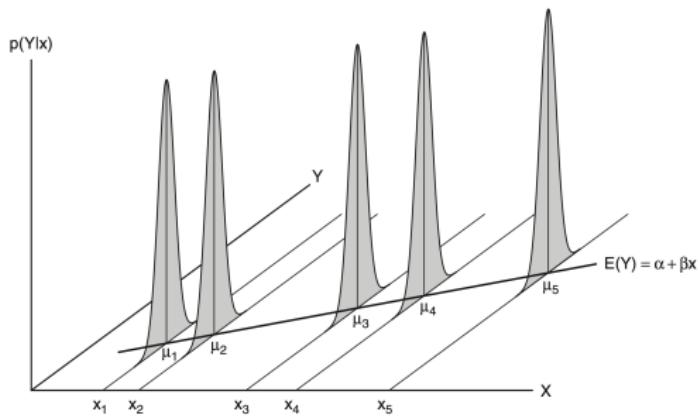
- In other words:

$$P(Y_i = y_i) = p(y_i|\theta_i, \phi),$$

$$\mu_i \equiv \mathbb{E}(Y_i) = a(\theta_i).$$

- Each point is described by factors $(X_1, X_2, \dots, X_n) \in \mathbb{R}^n$ - feature vectors.
- Model expected value μ_i for each point via linear relationship with link function g :

$$\mu = g^{-1}(\eta) = g^{-1}(w_1 X_1 + w_2 X_2 + \dots + w_n X_n).$$

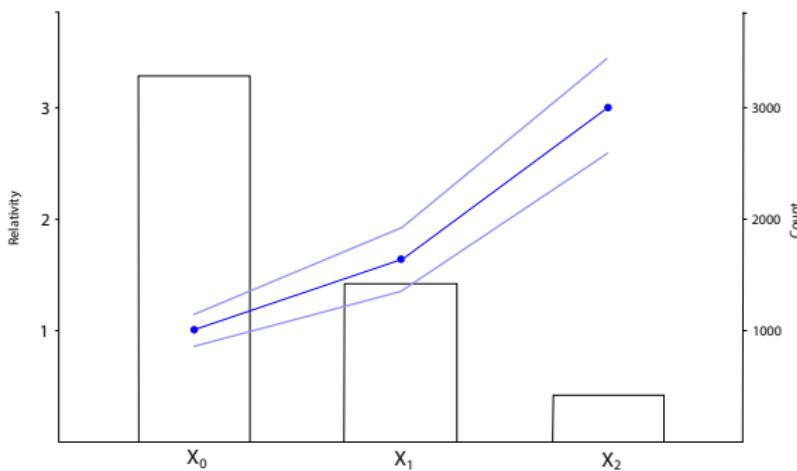


[source](#)

- The model can be (almost ☺) completely determined by parameters w_i .
- Estimate $\mathbf{w} = \{w_1, \dots, w_n\}$ via maximizing (log-) likelihood function:

$$\ell(\mathbf{w}) = \log p(\mathcal{D}|\mathbf{w}) = \log \prod_i p(y_i|\theta_i, \phi) = \sum_i \log p(y_i|\mathbf{x}_i, \theta_i, \phi),$$

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \ell(\mathbf{w}).$$



- Bernoulli distribution

$$p(y_i) = \text{Ber}(y_i|\mu) = \mu^y(1-\mu)^{1-y} = \begin{cases} 1 - \mu & \text{if } y=0 \\ \mu & \text{if } y=1 \end{cases}$$

- We want output to be in range from 0 to 1:

$$\mu \in [0, 1]$$

- Apply sigmoid (activation function) to satisfy this condition:

$$\mu = \sigma(\eta) = 1 / (1 + e^{-\eta}),$$

$$\eta = w_1 X_1 + w_2 X_2 + \dots + w_n X_n.$$

- Writing down log-likelihood for Bernoulli distribution leads to **LogLoss**:

$$\ell = \sum_i \log p(y_i | \mu_i) = - \sum_i (-y_i \log \mu_i - (1 - y_i) \log (1 - \mu_i)).$$

$$\text{LogLoss} = \sum_i (-y_i \log \mu_i - (1 - y_i) \log (1 - \mu_i)) \longrightarrow \text{minimize}$$

- LogLoss on 1 observation:

$$\text{LogLoss} = \begin{cases} -\log \mu_i & \text{if } y_i = 1, \\ -\log(1 - \mu_i) & \text{if } y_i = 0. \end{cases}$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}, \quad \mathbf{x}^T = [x_1 \ x_2 \ \dots \ x_n].$$

$$\boxed{\eta = \mathbf{w}^T \mathbf{x} \in \mathbb{R}}$$

$$(1 \times n) \cdot (n \times 1) = (1 \times 1) \in \mathbb{R}$$

$$\eta = [\mathbf{w}_1 \ \mathbf{w}_2 \ \dots \ \mathbf{w}_n] \cdot \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} = \mathbf{w}_1 x_1 + \mathbf{w}_2 x_2 + \dots + \mathbf{w}_n x_n.$$

$$\eta = \mathbf{X}\mathbf{w}$$

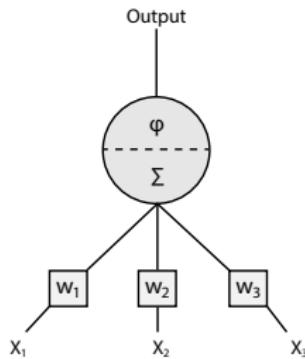
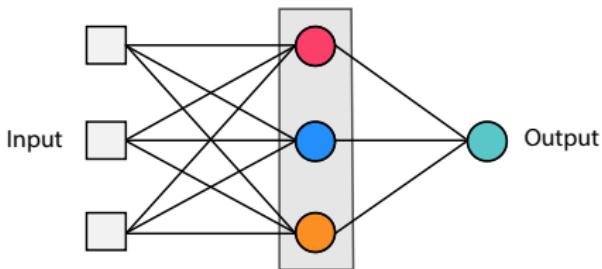
$$(N \times n) \cdot (n \times 1) = (N \times 1) \in \mathbb{R}^N$$

$$\begin{array}{c}
 \text{Feature vectors} \\
 \overbrace{\quad\quad\quad\quad\quad}^{\text{Feature vectors}} \quad\quad\quad\quad\quad \text{Weights} \\
 \overbrace{\quad\quad\quad\quad\quad}^{\text{Weights}}
 \end{array}$$

$$\begin{array}{l}
 \text{Obs}_1 \quad \left[\begin{array}{cccc} \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_n \end{array} \right] \cdot \left[\begin{array}{c} w_1 \\ w_2 \\ \dots \\ w_n \end{array} \right] \\
 \text{Obs}_2 \quad \left[\begin{array}{cccc} x_{11} & x_{12} & \dots & x_{1n} \\ \boxed{x_{21}} & x_{22} & \dots & x_{2n} \end{array} \right] \cdot \left[\begin{array}{c} w_1 \\ w_2 \\ \dots \\ w_n \end{array} \right] \\
 \dots \\
 \text{Obs}_N \quad \left[\begin{array}{cccc} x_{N1} & x_{N2} & \dots & x_{Nn} \end{array} \right] \cdot \left[\begin{array}{c} w_1 \\ w_2 \\ \dots \\ w_n \end{array} \right]
 \end{array}$$

$$\begin{array}{l}
 \text{Obs}_1 \quad \left[\begin{array}{c} w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n} \end{array} \right] \\
 \text{Obs}_2 \quad \left[\begin{array}{c} \boxed{w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n}} \\ \dots \end{array} \right] \\
 \text{Obs}_N \quad \left[\begin{array}{c} w_1 x_{N1} + w_2 x_{N2} + \dots + w_n x_{Nn} \end{array} \right]
 \end{array}$$

Linear Layer (Projection)



$$\text{Output} : \phi(w_1x_1 + w_2x_2 + w_3x_3) = \phi(\mathbf{x}^T \mathbf{w})$$

Neurons

$$\left[\begin{array}{cccc} N_1 & N_2 & \dots & N_d \\ w_{11} & w_{12} & \dots & w_{1d} \\ w_{21} & w_{22} & \dots & w_{2d} \\ \dots & \dots & \dots & \dots \\ w_{n1} & w_{n2} & \dots & w_{nd} \end{array} \right]$$

$$\eta = \mathbf{W}^T \mathbf{x} \in \mathbb{R}^d$$

$$(d \times n) \cdot (n \times 1) = (d \times 1) \in \mathbb{R}^d$$

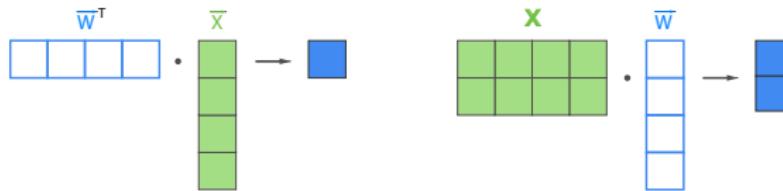
$$\mathcal{O} = \mathbf{X}\mathbf{W}$$

$$(N \times n) \cdot (n \times d) = (N \times d) \in \mathbb{R}^{N \times d}$$

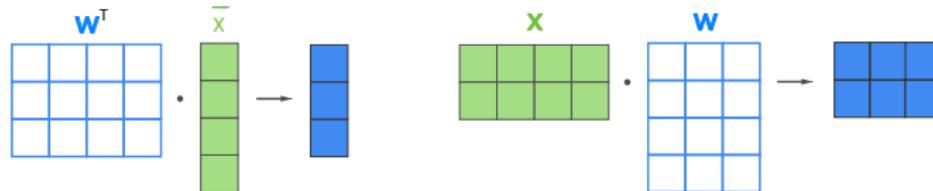
	Feature vectors				Neurons				
	X_1	X_2	\dots	X_n	N_1	N_2	\dots	N_d	
Obs ₁	x_{11}	x_{12}	\dots	x_{1n}	w_{11}	w_{12}	\dots	w_{1d}	w_1
Obs ₂	x_{21}	x_{22}	\dots	x_{2n}	w_{21}	w_{22}	\dots	w_{2d}	w_2
\dots	\dots	\dots	\dots	\dots	\dots	\dots	\dots	\dots	\dots
Obs _N	x_{N1}	x_{N2}	\dots	x_{Nn}	w_{n1}	w_{n2}	\dots	w_{nd}	w_n

	N_1	\dots	N_d	
Obs ₁	$x_{11}w_{11} + x_{12}w_{21} + \dots + x_{1n}w_{n1}$	\dots	$x_{11}w_{1d} + x_{12}w_{2d} + \dots + x_{1n}w_{nd}$	
Obs ₂	$x_{21}w_{11} + x_{22}w_{21} + \dots + x_{2n}w_{n1}$	\dots	$x_{21}w_{1d} + x_{22}w_{2d} + \dots + x_{2n}w_{nd}$	
\dots	\dots	\dots	\dots	
Obs _N	$x_{N1}w_{11} + x_{N2}w_{21} + \dots + x_{Nn}w_{n1}$	\dots	$x_{N1}w_{1d} + x_{N2}w_{2d} + \dots + x_{Nn}w_{nd}$	

LINEAR MODEL



LINEAR LAYER



PyTorch: Linear Layer

```
In 1 1 import torch
      2 import torch.nn as nn

In 2 1 linear = nn.Linear(4, 2, bias=False) # in features: 4, out_features: 2

In 3 1 input = torch.randn(3, 4)
      2 input

Out 3 tensor([[-0.2782,  0.0798, -0.3552, -0.2618],
             [ 0.6665,  1.9803,  0.2063, -1.0990],
             [-0.2645, -0.0047, -0.5869, -0.3289]])

In 4 1 # [3,4] * [2,4]^T = [3, 2]
      2 output = linear(input)
      3 output

Out 4 tensor([[ 0.1267, -0.0859],
             [ 0.8247,  0.8660],
             [ 0.0721, -0.0997]], grad_fn=<MmBackward>)

In 5 1 linear.weight

Out 5 Parameter containing:
tensor([[-0.3313,  0.4970,  0.0487, -0.0466],
       [ 0.4462,  0.1334,  0.1113, -0.2562]], requires_grad=True)
```

```
In 1 1 import torch
      2 import torch.nn as nn

In 2 1 class Model(nn.Module):
      2     """ Simple model: 2 linear layers and ReLU activation. """
      3     def __init__(self,
      4             in_dim,
      5             hid_dim):
      6
      7         super().__init__()
      8
      9         self.in_dim = in_dim
     10        self.linear1 = nn.Linear(in_dim, hid_dim, bias=False)
     11        self.linear2 = nn.Linear(hid_dim, 1, bias=False)
     12        self.relu = nn.ReLU()
     13
     14    def forward(self, input):
     15        out = self.linear1(input) # [batch_size, hid_dim]
     16        out = self.relu(out) # activation function
     17        out = self.linear2(out) # [batch_size, 1]
     18
     19        return out
```

```
In 3 1 model = Model(in_dim=4,
2           hid_dim=2)

In 4 1 model.linear1.weight

Out 4      Parameter containing:
tensor([[ -0.1320,   0.0753,  -0.3863,   0.4888],
       [-0.2059,  -0.2908,   0.4562,  -0.0320]], requires_grad=True)

In 5 1 model.linear2.weight

Out 5      Parameter containing:
tensor([[ 0.3291,   0.3062]], requires_grad=True)

In 6 1 input = torch.randn(2,4, dtype=torch.float32)
2 input

Out 6      tensor([[-0.6946, -1.0389, -0.3978,  1.0273],
       [-0.0191,   0.8283,   0.3426,  2.6409]])

In 7 1 output = model(input)
2 output

Out 7      tensor([[ 0.2909,
       [ 0.4026]], grad_fn=<MmBackward>)
```

- Parameters are *Tensor* subclasses.
- They're automatically added to list of *Module* parameters when assigned as attributes in constructor.

```
In 2 1 class Model(nn.Module):
2
3     def __init__(self):
4         super().__init__()
5         self.W = nn.Parameter(torch.Tensor(2,2))
6         self.T = torch.zeros(2, 2, requires_grad=True)

In 3 1 model = Model()
2 for p in model.named_parameters():
3     print(p)

('W', Parameter containing:
 tensor([[0., 0.],
        [0., 0.]], requires_grad=True))

In 4 1 # initialize weights
2 for param in model.parameters():
3     stdv = 1. / math.sqrt(param.size(-1))
4     param.data.uniform_(-stdv, stdv)

In 5 1 import torch.optim as optim
2
3 # optimizer will update only registered Parameters of the model
4 optimizer = optim.Adam(model.parameters())
```

Data	Label Encoding	OHE vectors
Color	Color	Blue Red Grey
Blue	1	1 0 0
Red	2	0 1 0
Grey	3	0 0 1
Grey	3	0 0 1
Blue	1	1 0 0

$$\mathbf{x} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \mathbf{W} = \begin{bmatrix} 5 & 2 & 1 & 0 \\ 3 & 6 & 9 & 5 \\ 2 & 1 & 0 & 7 \end{bmatrix}.$$

$$\mathbf{W}^T \mathbf{x} = ?$$

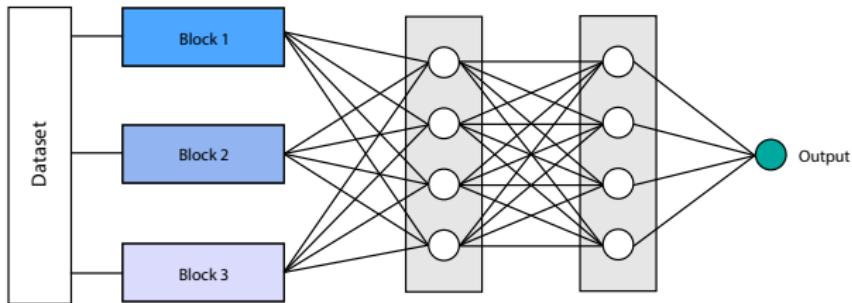
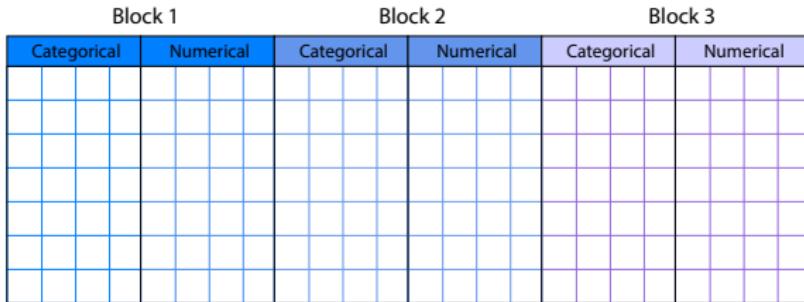
$$\mathbf{W}^T \mathbf{x} = \begin{bmatrix} 5 & 3 & 2 \\ 2 & 6 & 1 \\ 1 & 9 & 0 \\ 0 & 5 & 7 \end{bmatrix} \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 3 \\ 6 \\ 9 \\ 5 \end{bmatrix}.$$

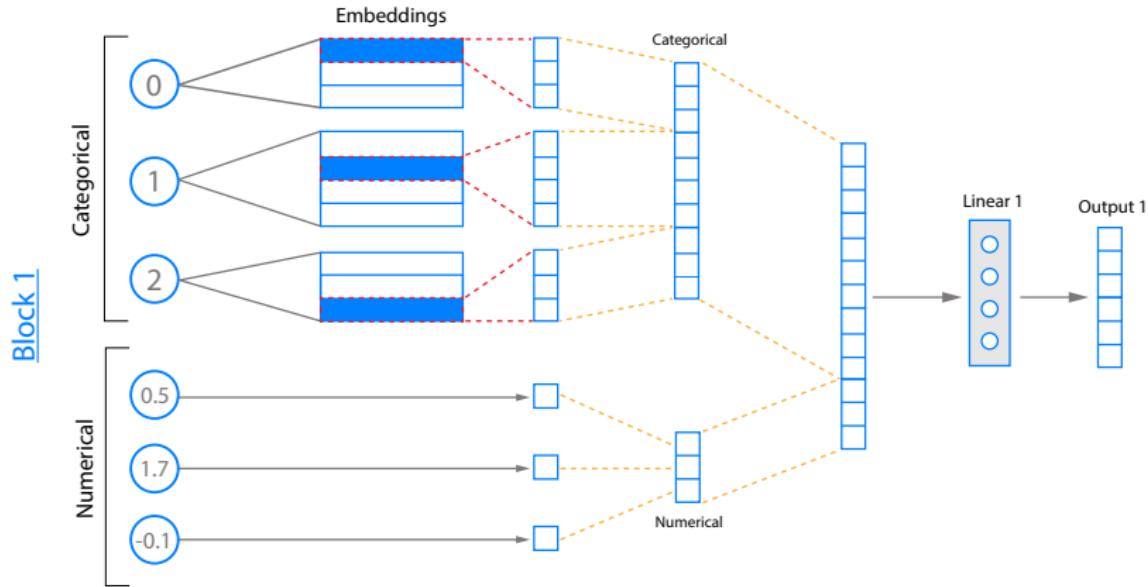
$$\mathbf{W} = \begin{bmatrix} 5 & 2 & 1 & 0 \\ 3 & 6 & 9 & 5 \\ 2 & 1 & 0 & 7 \end{bmatrix}.$$

Good news: we don't need to OHE categorical variables.

Label-encode categories + initialize embedding matrix \rightarrow learnable lookup table.

Categorical Factors





Concatenate outputs from all blocks and pass to the remainder layers of NN.

Entity Embeddings of Categorical Variables
[arXiv:1604.06737 \[cs.LG\]](https://arxiv.org/abs/1604.06737)

Embedding layer in PyTorch:

```
In 1 1 import torch
      2 import torch.nn as nn

In 2 1 # an Embedding layer: 4 vectors of size 3
      2 embedding = nn.Embedding(4, 3)
      3 embedding.weight

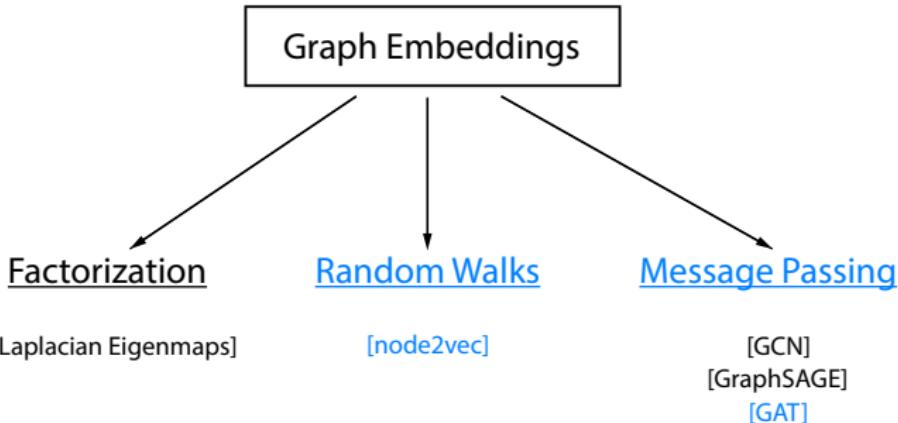
Out 2 Parameter containing:
tensor([[ 1.0122,  2.0579, -0.0501],
       [ 3.3124, -0.9464,  1.3357],
       [-0.2790,  0.2287, -2.1280],
       [ 0.7080, -1.1197, -0.1263]], requires_grad=True)

In 3 1 # a batch of 2 samples of 3 indices each
      2 input = torch.LongTensor([[0,1,2],[3,2,1]])
      3 embedding(input)

Out 3 tensor([[[ 1.0122,  2.0579, -0.0501],
              [ 3.3124, -0.9464,  1.3357],
              [-0.2790,  0.2287, -2.1280]],

             [[ 0.7080, -1.1197, -0.1263],
              [-0.2790,  0.2287, -2.1280],
              [ 3.3124, -0.9464,  1.3357]]], grad_fn=<EmbeddingBackward>)
```

Graph embeddings



Representation Learning on Graphs: Methods and Applications

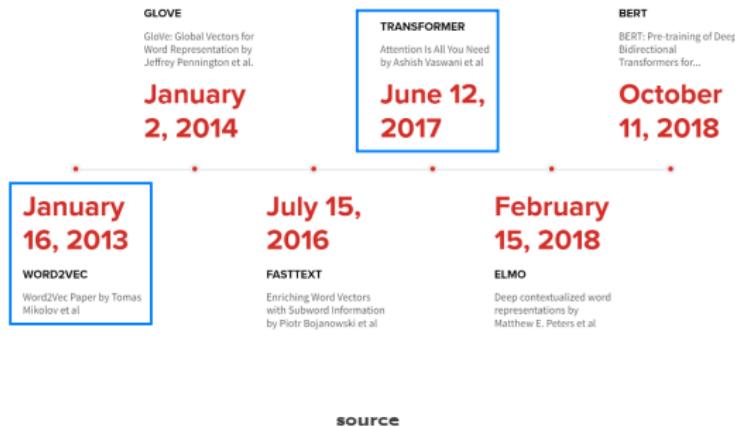
[arXiv:1709.05584 \[cs.SI\]](https://arxiv.org/abs/1709.05584)

Graph Representation Learning Book [Hamilton]

https://www.cs.mcgill.ca/~wlh/grl_book/



- Many approaches in modern Deep Learning are closely related to ideas from Natural Language Processing (**NLP**).
- Timeline of some major NLP methods:



[source](#)

- Check the [Stanford CS224N](#) course.

- Distributional Hypothesis: words that occur in the same contexts tend to have similar meanings [Harris, 1954]
- "A word is characterized by the company it keeps" [Firth, 1957]
- Can you guess the meaning of the word? [source](#)
 - ① A bottle of **tezgüino** is on the table.
 - ② Everyone likes **tezgüino**.
 - ③ **Tezgüino** makes you drunk.
 - ④ We make **tezgüino** for fun.

word2vec

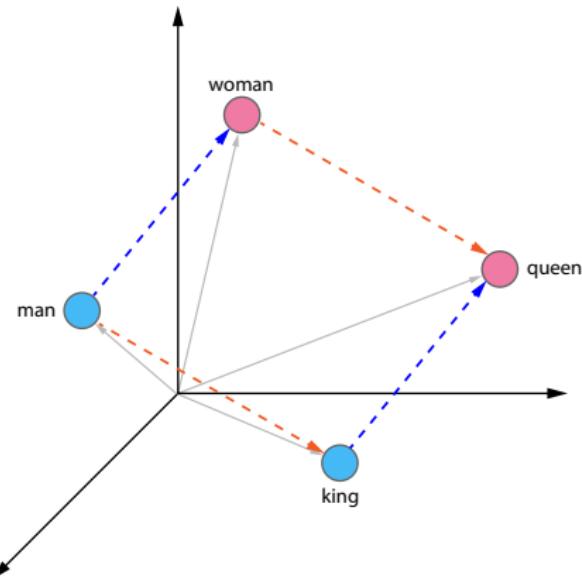
Efficient Estimation of Word Representations in Vector Space

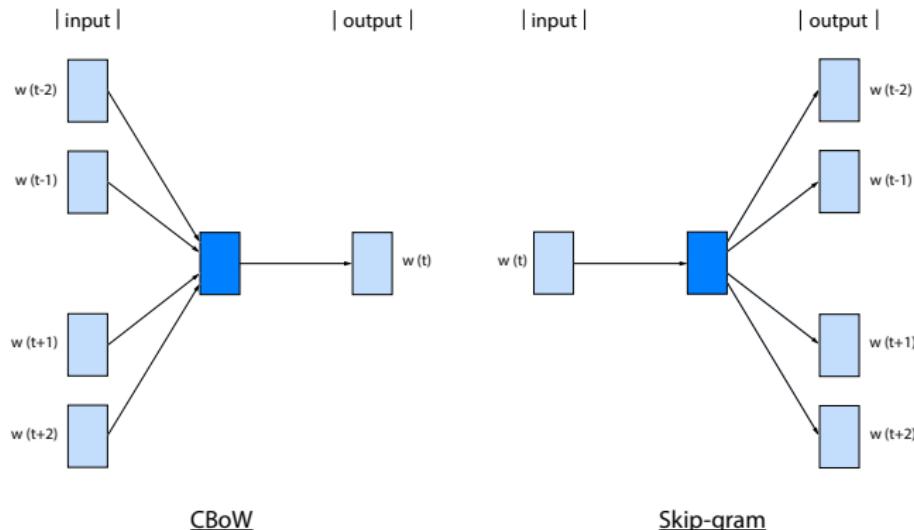
[arXiv:1301.3781 \[cs.CL\]](https://arxiv.org/abs/1301.3781)

word2vec Parameter Learning Explained

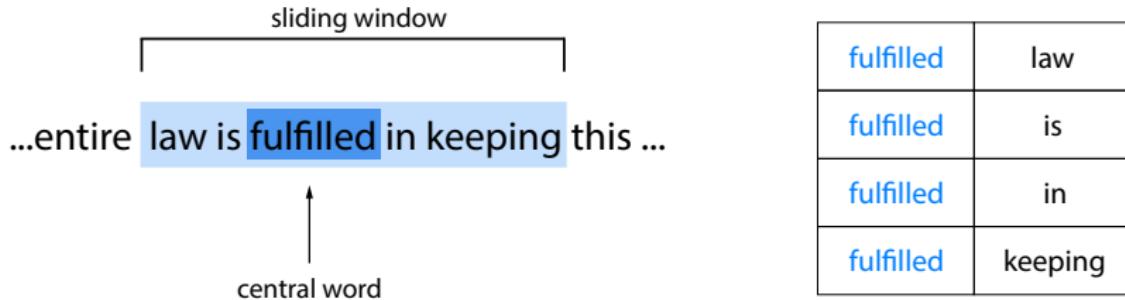
[arXiv:1411.2738 \[cs.CL\]](https://arxiv.org/abs/1411.2738)

$$\text{vector[King]} - \text{vector[Man]} + \text{vector[Woman]} = \text{vector[Queen]}$$



CBoWSkip-gram

For the entire law is fulfilled in keeping this one command:
"Love your neighbor as yourself".

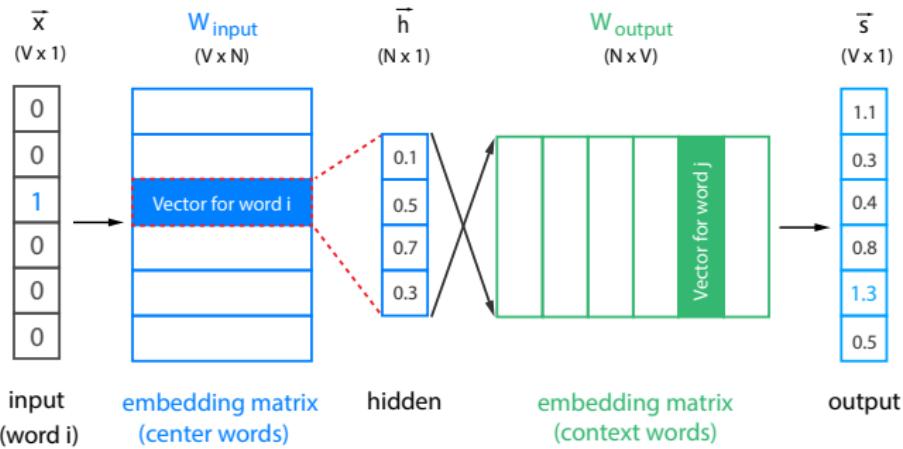


Keep on sliding

...entire law is fulfilled in keeping this ...

... law is fulfilled in keeping this one ...

fulfilled	law
fulfilled	is
fulfilled	in
fulfilled	keeping
in	is
in	fulfilled
in	keeping
in	this



$$\mathbf{h} = \mathbf{W}_{\text{input}}^T \mathbf{x} = \mathbf{v}_{w_i}^T.$$

$$\mathbf{s} = \mathbf{W}_{\text{output}}^T \mathbf{h}.$$

$$s_j = \langle \mathbf{v}_{w_j}'; \mathbf{v}_{w_i} \rangle$$

Similarity of central and context words embeddings

Softmax function / Boltzmann distribution

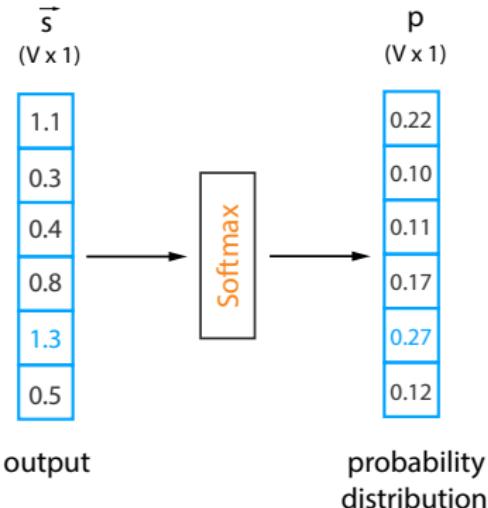
$$\mathbf{z} = (z_1, \dots, z_K) \in \mathbb{R}^K.$$

$$\sigma(\mathbf{z})_i = \frac{e^{\beta z_i}}{\sum_{j=1}^K e^{\beta z_j}}, \quad \beta = 1/T.$$

- $\sigma(\mathbf{z}_i) : \mathbb{R}^K \rightarrow [0, 1]^K$.
- $\sum_{i=1}^K \sigma(\mathbf{z})_i = 1$.
- Normalizes output to **probability distribution** over predicted classes.
- T - temperature.

At **low temperatures** the distribution puts most of its probability mass in the most probable state.

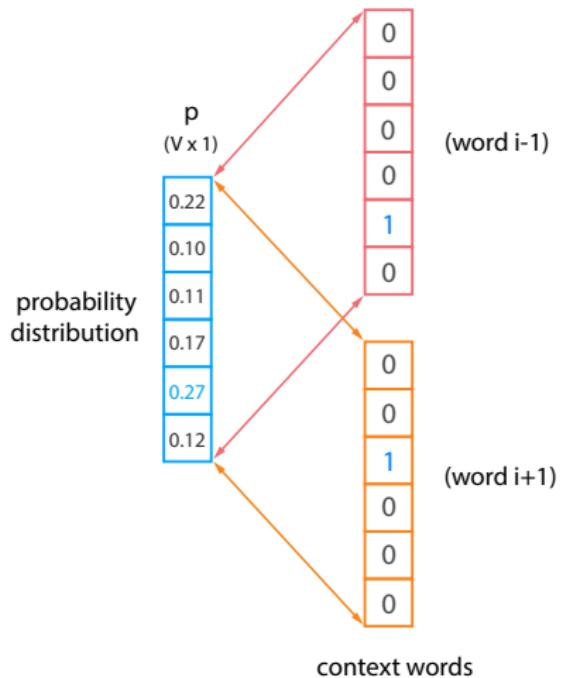
At **high temperatures** it spreads the mass more uniformly.



- Probability of observing the actual output word w_O given the input word w_I :

$$p(w_O|w_I) = \frac{\exp(\langle v'_{w_O}; v_{w_I} \rangle)}{\sum_{j=1}^V \exp(\langle v'_{w_j}; v_{w_I} \rangle)}.$$

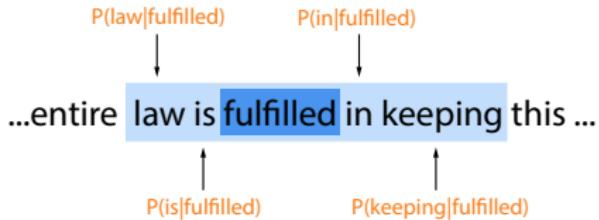
- V can be extremely large. Calculating denominator is computationally complex.



- Compute probabilities of a context words for each central word in corpus:

$$\mathcal{L}_\theta = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t).$$

$t = 1, \dots, T$ – positions of words,
 m – sliding window size.

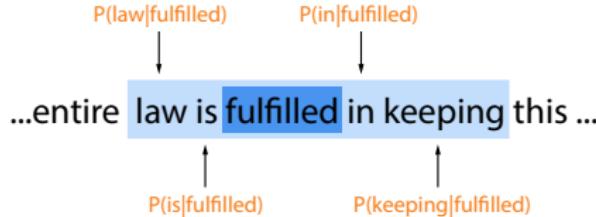


- Loss function $J(\theta)$ - average negative log-likelihood:

$$J(\theta) \triangleq -\frac{1}{T} \log \mathcal{L}_\theta = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_t).$$

- $w_t = \text{fulfilled}$, $w_{t+j} = \text{law}$, $j = -2$:

$$-\log P(\text{law} | \text{fulfilled}) = -\log \frac{\exp (\langle \mathbf{v}'_{\text{law}}; \mathbf{v}_{\text{fulfilled}} \rangle)}{\sum_{j=1}^V \exp (\langle \mathbf{v}'_j; \mathbf{v}_{\text{fulfilled}} \rangle)}.$$



- Loss term:

$$J_{t,j}(\theta) = -\langle \mathbf{v}'_{\text{law}}; \mathbf{v}_{\text{fulfilled}} \rangle + \log \sum_{j=1}^V \exp (\langle \mathbf{v}'_j; \mathbf{v}_{\text{fulfilled}} \rangle).$$

- How to minimize Loss:

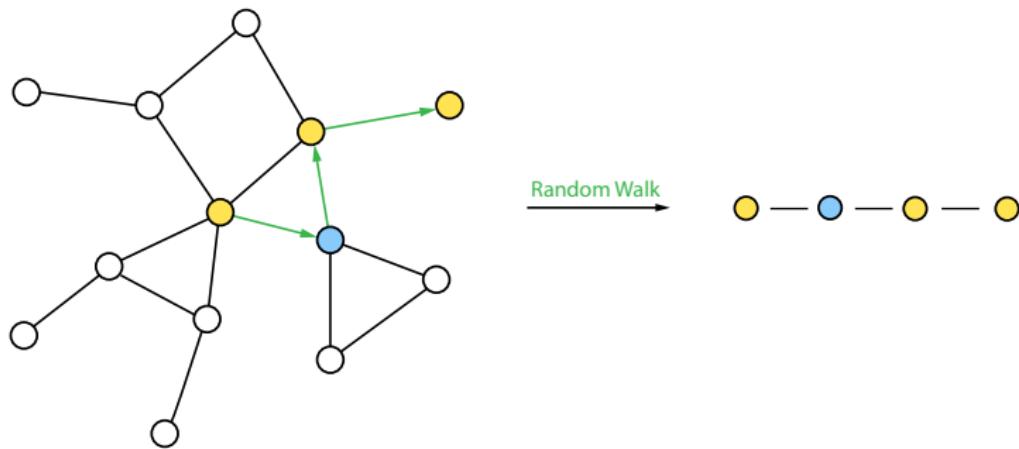
- increase similarity between central and context vector: $\langle \mathbf{v}'_{\text{law}}; \mathbf{v}_{\text{fulfilled}} \rangle$.
- decrease similarity between central and non-context words: $\log \sum_j \exp(\dots)$.

node2vec

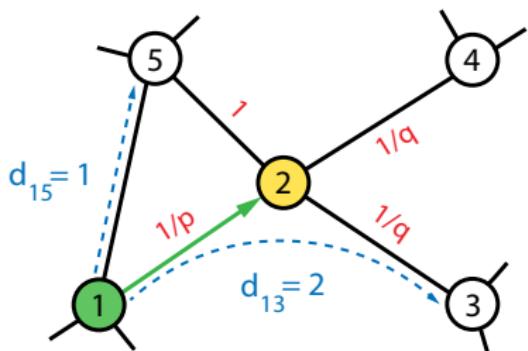
node2vec: Scalable Feature Learning for Networks

[arXiv:1607.00653 \[cs.SI\]](https://arxiv.org/abs/1607.00653)

- Unsupervised / self-supervised embeddings - task independent.
 - ① Not using node labels.
 - ② Not using node features.
- Source node $u \rightarrow$ neighbourhood $N_S(u)$:



- Biased random walk of fixed length: interpolate between local (BFS) and global (DFS) views of neighbourhoods.
- Parameters:
 - ➊ p - return parameter
 - ➋ q - in-out parameter

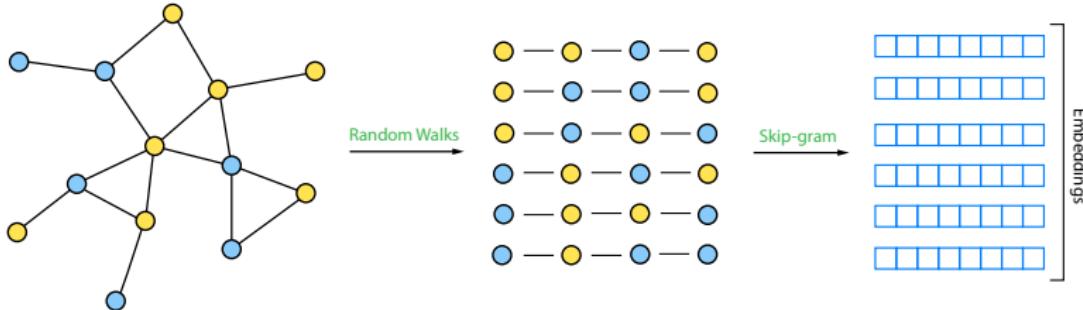


Transitioned $1 \rightarrow 2$,
evaluating next step.

Unnormalized transition
probabilities:

$$\alpha(t, x) = \begin{cases} \frac{1}{p} & \text{if } d_{tx} = 0 \\ 1 & \text{if } d_{tx} = 1 \\ \frac{1}{q} & \text{if } d_{tx} = 2 \end{cases}$$

d_{tx} - shortest distance between
 t and x .



Maximize log-probability of observing neighbourhood $N_S(u)$ for node u :

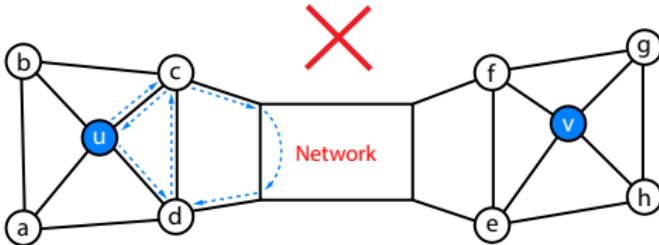
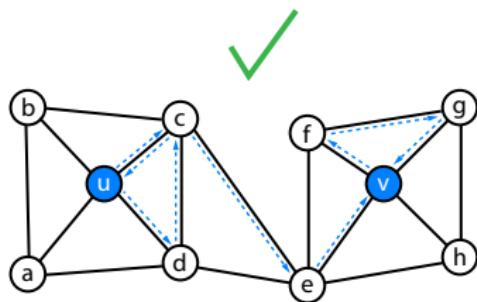
$$J(\theta) \sim -\log \mathcal{L}_\theta = -\sum_{u \in V} \log P(N_S(u) | \mathbf{z}_u) = -\sum_{u \in V} \sum_{v \in N_S(u)} \log P(v | \mathbf{z}_u).$$

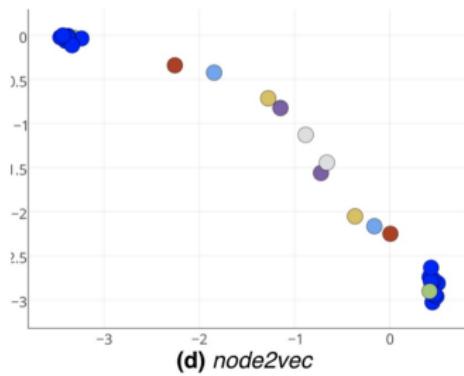
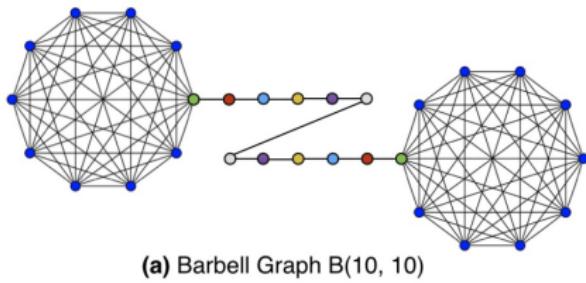
$$P(N_S(u) | \mathbf{z}_u) = \prod_{n_i \in N_S(u)} P(n_i | \mathbf{z}_u)$$

Probability parametrized using softmax:

$$P(v | \mathbf{z}_u) = \frac{\exp(\langle \mathbf{z}_v; \mathbf{z}_u \rangle)}{\sum_{n \in V} \exp(\langle \mathbf{z}_n; \mathbf{z}_u \rangle)}$$

- node2vec tend to fail in structural equivalence tasks.
- Nodes that are far in the network tend to be separated in latent representation.





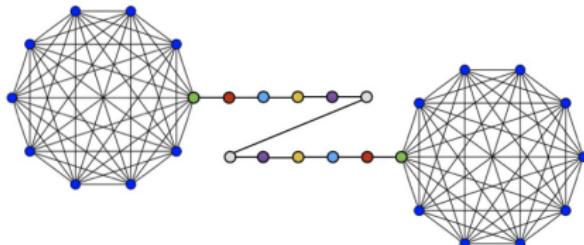
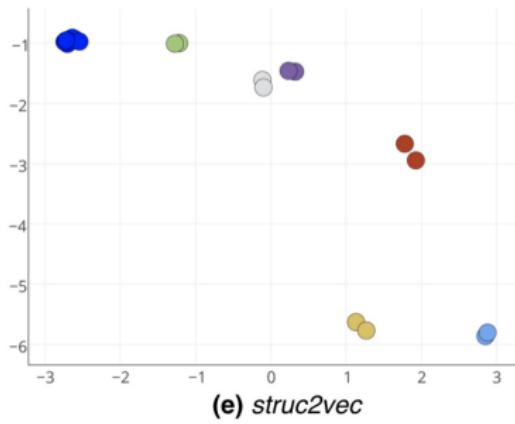
struc2vec

struc2vec: Learning Node Representations from Structural Identity

[arXiv:1704.03165 \[cs.SI\]](https://arxiv.org/abs/1704.03165)

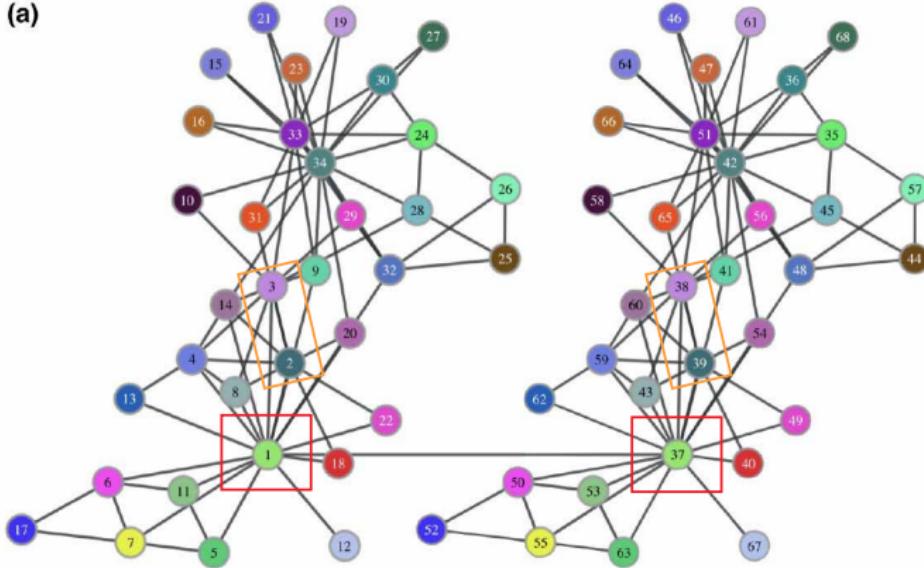
<https://github.com/leoribeiro/struc2vec>

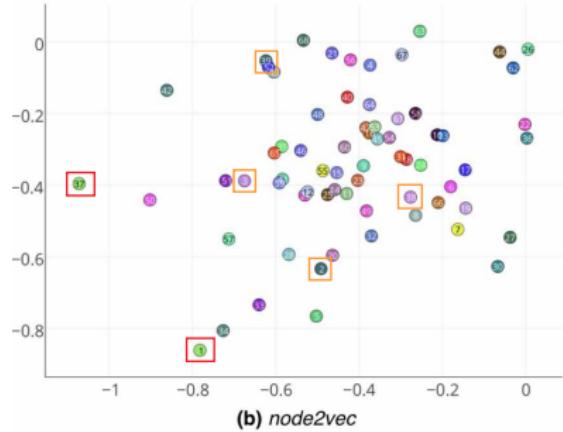
- Latent representations for the structural identity of nodes
- Biased random walks on multilayer graph

(a) Barbell Graph $B(10, 10)$ 

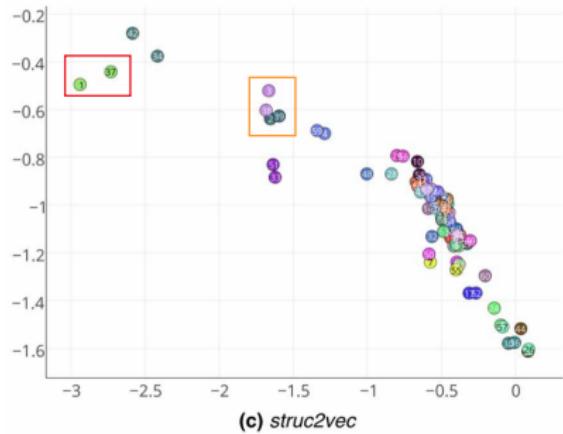
(e) struc2vec

(a)





(b) node2vec



(c) struc2vec

Message Passing

Neural Message Passing for Quantum Chemistry

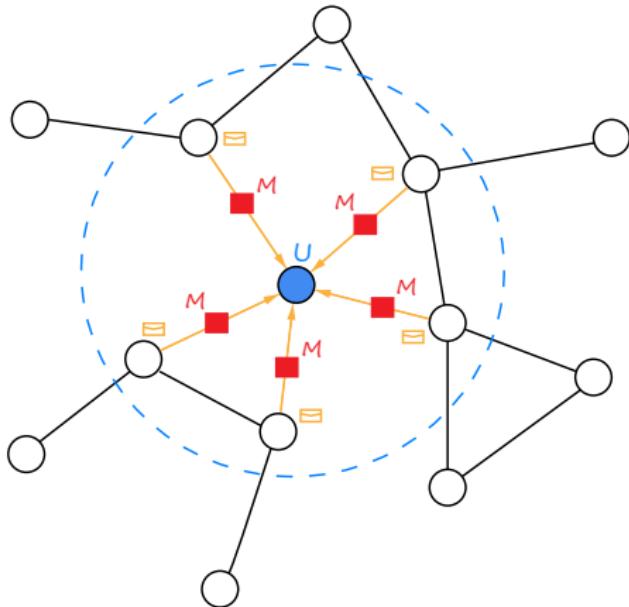
[arXiv:1704.01212 \[cs.LG\]](https://arxiv.org/abs/1704.01212)

Embedding $h_i^{(t)}$ for each node V_i is obtained by aggregating a function \mathcal{M} of features $h_j^{(t-1)}$ (and possibly edges) of all neighbouring nodes V_j with the node V_i itself:

$$h_i^{(t)} = \mathcal{U}^{(t)} \left(h_i^{(t-1)}, \square_{j \in \mathcal{N}(i)} \mathcal{M}^{(t)}(h_i^{(t-1)}, h_j^{(t-1)}, e_{\{j,i\}}) \right),$$

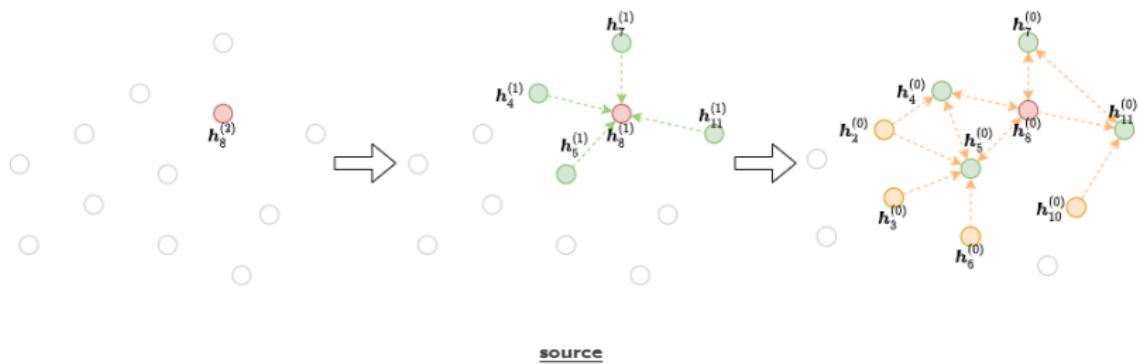
\square - aggregation function (ex. \sum_j),
 \mathcal{U}, \mathcal{M} - update/message differentiable functions (ex. MLP),
 $e_{\{j,i\}}$ - edge features from node j to node i .

$$\underline{h_i^{(t)} = \mathcal{U}^{(t)} \left(h_i^{(t-1)}, \square_{j \in \mathcal{N}(i)} \mathcal{M}^{(t)}(h_i^{(t-1)}, h_j^{(t-1)}, e_{\{j,i\}}) \right)}$$



Iterative process for 2-layer message passing (2-hop neighbourhood):

At each layer we obtain a bipartite graph structure



Graph Neural Networks: A Review of Methods and Applications [arXiv:1812.08434 \[cs.LG\]](https://arxiv.org/abs/1812.08434)

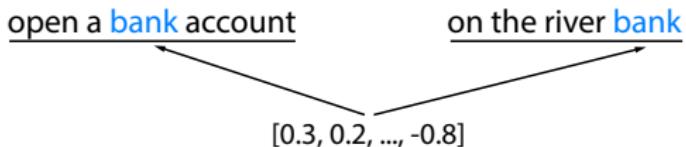
Name	Variant	Aggregator	Updater
Spectral Methods	ChebNet	$\mathbf{N}_k = \mathbf{T}_k(\tilde{\mathbf{L}})\mathbf{X}$	$\mathbf{H} = \sum_{k=0}^K \mathbf{N}_k \Theta_k$
	1 st -order model	$\mathbf{N}_0 = \mathbf{X}$ $\mathbf{N}_1 = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = \mathbf{N}_0 \Theta_0 + \mathbf{N}_1 \Theta_1$
	Single parameter	$\mathbf{N} = (\mathbf{I}_N + \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}) \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$
Non-spectral Methods	GCN	$\mathbf{N} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X}$	$\mathbf{H} = \mathbf{N} \Theta$ Graph Convolutions
	Neural FPs	$\mathbf{h}'_{\mathcal{N}_v} = \mathbf{h}_v^{t-1} + \sum_{k=1}^{\mathcal{N}_v} \mathbf{h}_k^{t-1}$	$\mathbf{h}_v^t = \sigma(\mathbf{h}'_{\mathcal{N}_v} \mathbf{W}_L^{\mathcal{N}_v})$
	DCNN	Node classification: $\mathbf{N} = \mathbf{P}^* \mathbf{X}$ Graph classification: $\mathbf{N} = \mathbf{1}_N^T \mathbf{P}^* \mathbf{X} / N$	$\mathbf{H} = f(\mathbf{W}^c \odot \mathbf{N})$
Graph Attention Networks	GraphSAGE	$\mathbf{h}_{\mathcal{N}_v}^t = \text{AGGREGATE}_t(\{\mathbf{h}_u^{t-1}, \forall u \in \mathcal{N}_v\})$	$\mathbf{h}_v^t = \sigma(\mathbf{W}^t \cdot [\mathbf{h}_v^{t-1} \ \mathbf{h}_{\mathcal{N}_v}^t])$
	GAT	$\alpha_{vk} = \frac{\exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_v \ \mathbf{W}\mathbf{h}_k]))}{\sum_{j \in \mathcal{N}_v} \exp(\text{LeakyReLU}(\mathbf{a}^T [\mathbf{W}\mathbf{h}_v \ \mathbf{W}\mathbf{h}_j]))}$ $\mathbf{h}_{\mathcal{N}_v}^t = \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{vk} \mathbf{W}\mathbf{h}_k)$ Multi-head concatenation: $\mathbf{h}_{\mathcal{N}_v}^t = \parallel_{m=1}^M \sigma(\sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k)$ Multi-head average: $\mathbf{h}_{\mathcal{N}_v}^t = \sigma\left(\frac{1}{M} \sum_{m=1}^M \sum_{k \in \mathcal{N}_v} \alpha_{vk}^m \mathbf{W}^m \mathbf{h}_k\right)$	$\mathbf{h}_v^t = \mathbf{h}_{\mathcal{N}_v}^t$ Attention mechanism
Gated Graph Neural Networks	GGNN	$\mathbf{h}_{\mathcal{N}_v}^t = \sum_{k \in \mathcal{N}_v} \mathbf{h}_k^{t-1} + \mathbf{b}$	$\mathbf{z}_v^t = \sigma(\mathbf{W}^z \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^z \mathbf{h}_v^{t-1})$ $\mathbf{r}_v^t = \sigma(\mathbf{W}^r \mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}^r \mathbf{h}_v^{t-1})$ $\hat{\mathbf{h}}_v^t = \tanh(\mathbf{W}\mathbf{h}_{\mathcal{N}_v}^t + \mathbf{U}(\mathbf{r}_v^t \odot \mathbf{h}_v^{t-1}))$ $\mathbf{h}_v^t = (1 - \mathbf{z}_v^t) \odot \mathbf{h}_v^{t-1} + \mathbf{z}_v^t \odot \hat{\mathbf{h}}_v^t$

Transformers

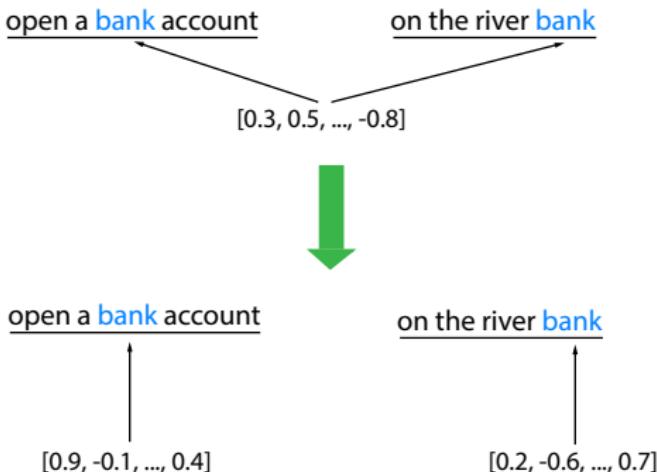
Attention Is All You Need

[arXiv:1706.03762 \[cs.CL\]](https://arxiv.org/abs/1706.03762)

- Word embeddings (word2vec, Glove) are *pre-trained* on text corpus based on co-occurrence statistics.
- Word embeddings are applied in a **context free** manner.

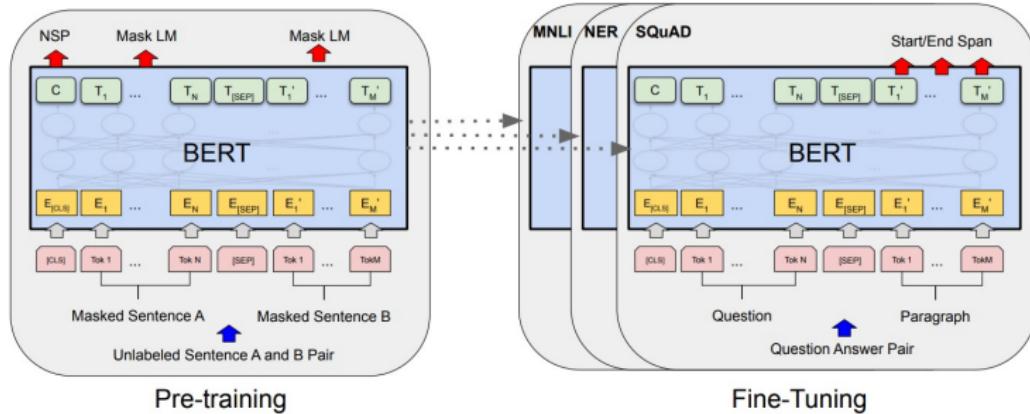


We need contextual embeddings



Bidirectional Encoder Representations from Transformers

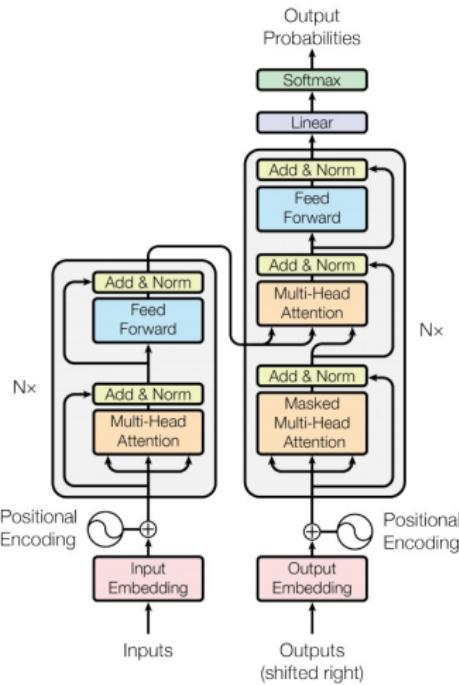
BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding
[arXiv:1810.04805 \[cs.CL\]](https://arxiv.org/abs/1810.04805)

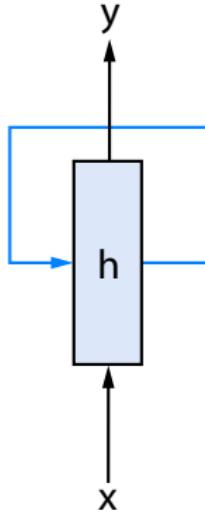


BERT pretraining - two unsupervised tasks:

- ❶ Masked LM: mask some input words and predict them.
- ❷ Next Sentence Prediction: predict whether *Sentence B* is actual sentence that proceeds *Sentence A*.

Transformer architecture





➊ Simplest RNN ($y_t = h_t$):

$$\mathbf{Y}_{(t)} = \phi(\mathbf{X}_{(t)}\mathbf{W} + \mathbf{Y}_{(t-1)}\mathbf{W}_y + \mathbf{b}).$$

➋ Rewrite in concatenated form:

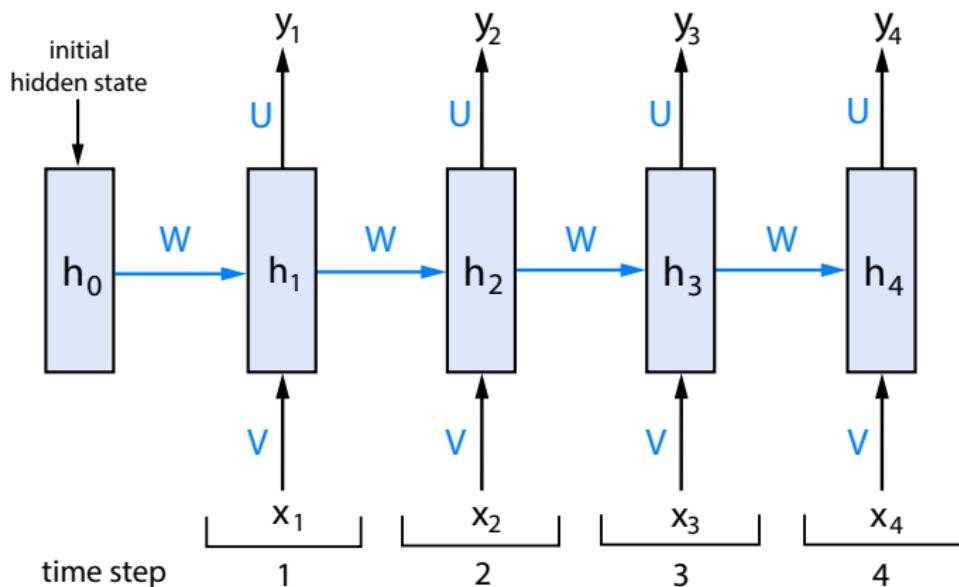
$$\mathbf{Y}_{(t)} = \phi([\mathbf{X}_{(t)} \ \mathbf{Y}_{(t-1)}]\mathbf{W} + \mathbf{b}),$$

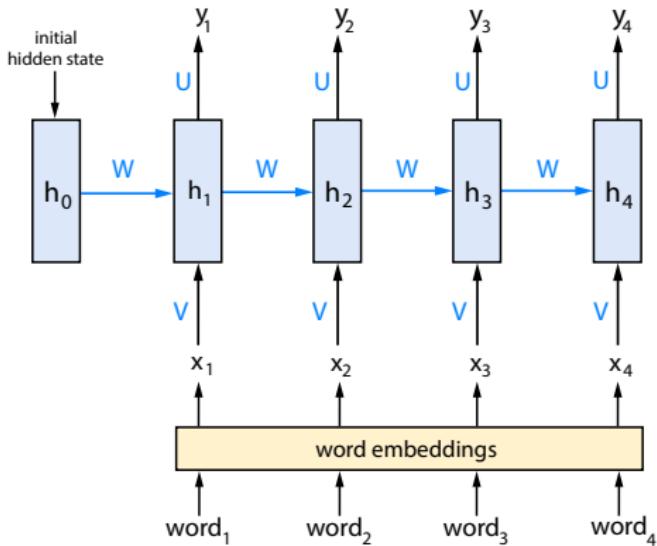
$$\mathbf{W} = \begin{bmatrix} \mathbf{W}_x \\ \mathbf{W}_y \end{bmatrix}.$$

➌ More complex:

- $y_t = g(h_t)$
- LSTM
- GRU
- Bidirectional
- Multilayer

Unfold the RNN!





Use cases:

- Tagging (ex: named entity recognition, part-of-speech tagging)
- Sentence classification (ex: sentiment analysis)

- Source sequence: $S = (x_1, \dots, x_S)$.
- Target sequence: $T = (y_1, \dots, y_T)$.
- Maximize log probability of correct translation T given source S in set \mathcal{S} :

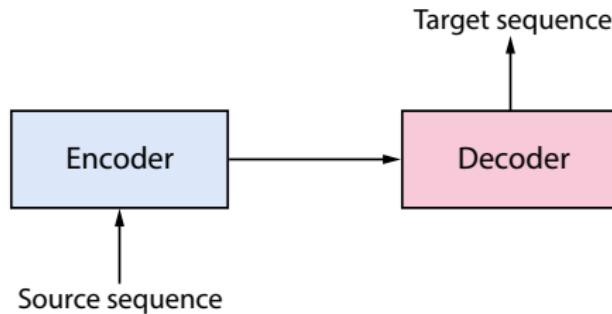
$$\mathcal{L} = \frac{1}{|\mathcal{S}|} \sum_{(T,S) \in \mathcal{S}} \log p(T | S)$$

- Neural network defines a probability over translation by decomposing the joint probability into the ordered conditionals:

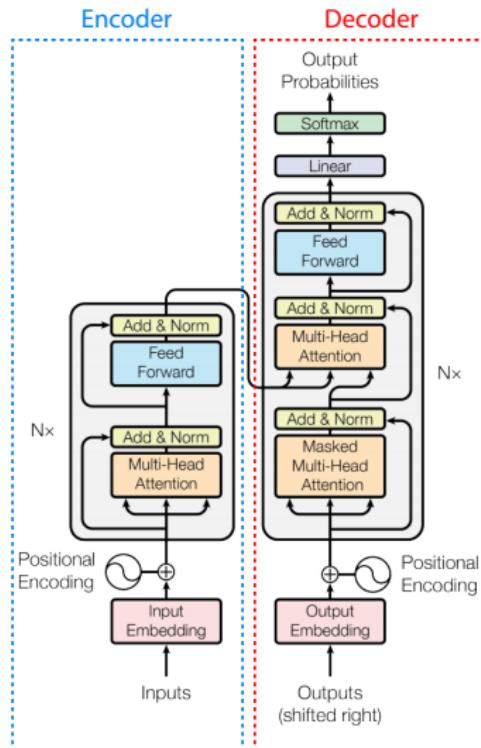
$$p(y_1, \dots, y_T | x_1, \dots, x_S) = \prod_{t=1}^T p(y_t | y_1, \dots, y_{t-1}, v),$$

v - representation of (x_1, \dots, x_S) .

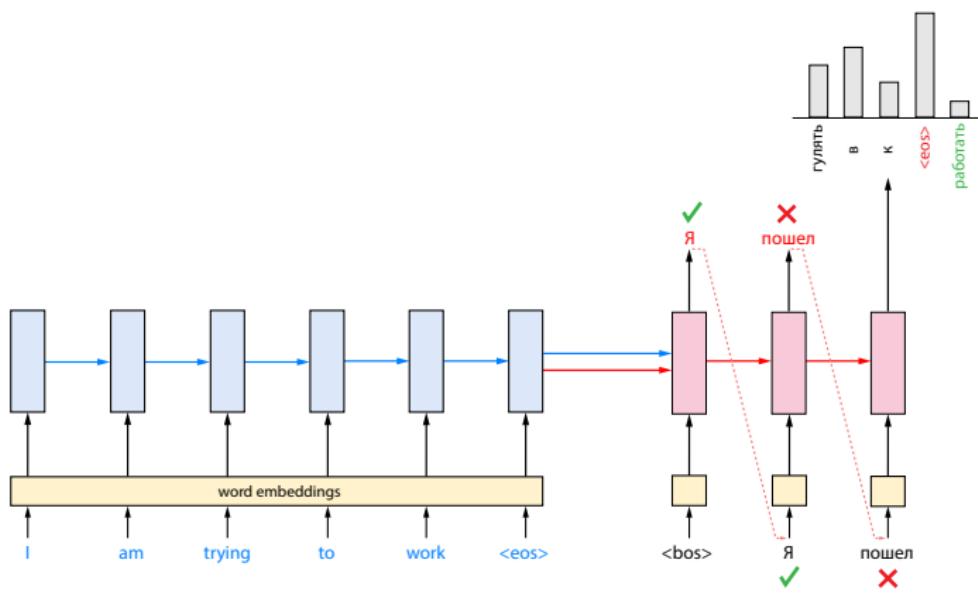
- **Encoder** produces a representation of source sequence and passes it to decoder.
- **Decoder** generates target sequence based on encoding.



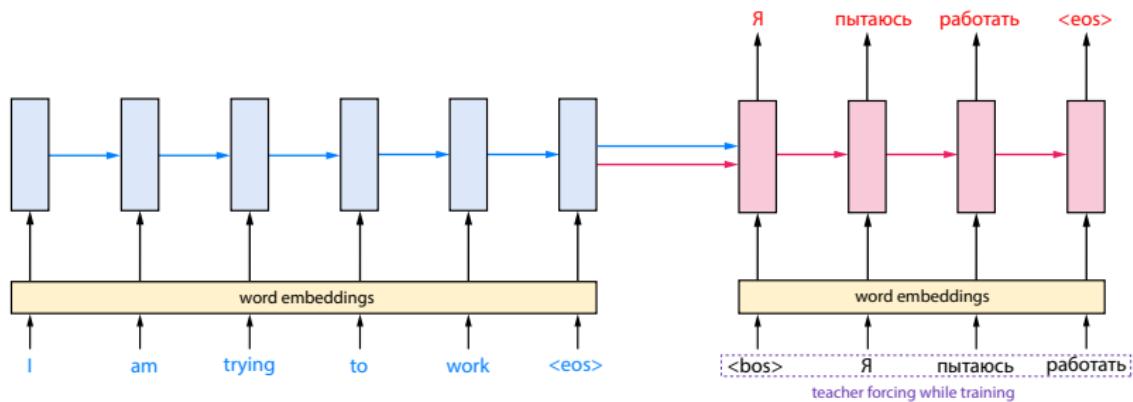
Transformer: Encoder-Decoder



- Simple case: two RNNs.
- Encoder's final state h_S - representation for the entire source sentence - serves as initial hidden state for Decoder RNN.



- Teacher forcing (training): replace the output $y(t)$ generated by network with the actual (expected) value in the training dataset.
- Bottleneck: all source information is encoded in single vector.



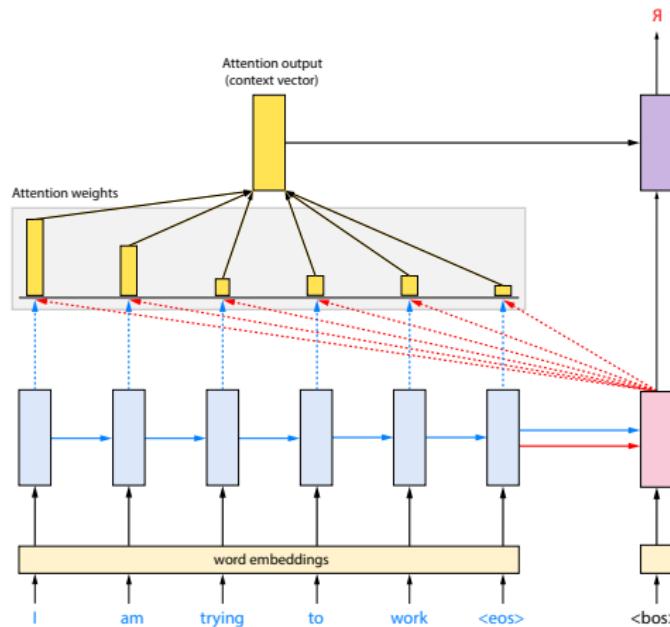
Attention

Neural Machine Translation by Jointly Learning to Align and Translate
[arXiv:1409.0473 \[cs.CL\]](#)

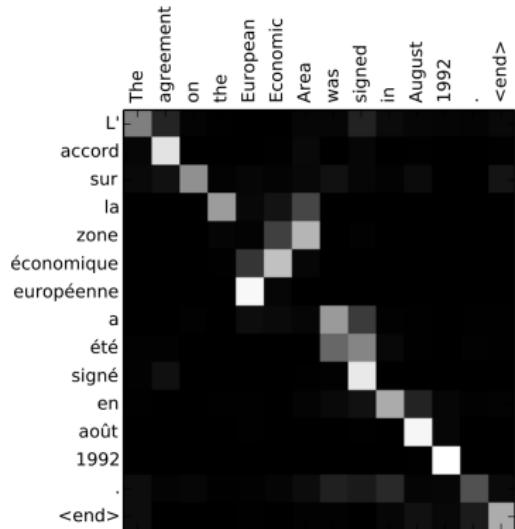
Effective Approaches to Attention-based Neural Machine Translation
[arXiv:1508.04025 \[cs.CL\]](#)

Attention Is All You Need
[arXiv:1706.03762 \[cs.CL\]](#)

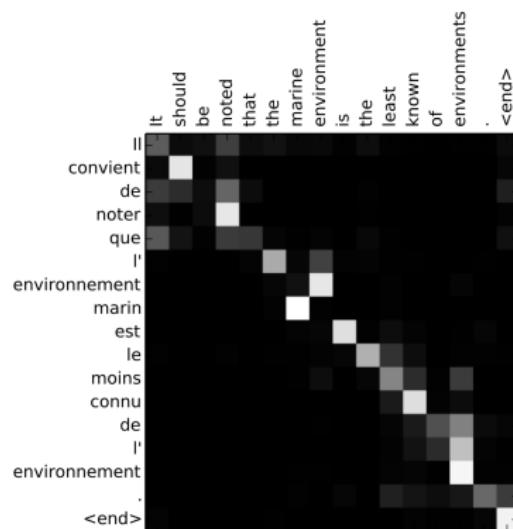
- At different decoder steps, Attention lets a model to focus on different parts of the input (according to their relevance for the current decoder state).



English → French



(a)



(b)

source

Document Classification

4 \longleftrightarrow star 5, 0 \longleftrightarrow star 1.

GT: 4 Prediction: 4

pork belly = delicious .

scallops ?

i do n't .

even .

like .

scallops , and these were a-m-a-z-i-n-g .

fun and tasty cocktails .

next time i 'm in phoenix , i will go

back here .

highly recommend .

GT: 0 Prediction: 0

terrible value .

ordered pasta entree .

.

\$ 16.95 good taste but size was an appetizer size .

.

no salad , no bread no vegetable .

this was .

our and tasty cocktails .

our second visit .

i will not go back .

Hierarchical Attention Networks for Document Classification

Calculation for 1 pair ($S \leftrightarrow T$) in batch

- d - hidden dimension of the model, m, n - seq lens.
- A set of m **query** vectors (come from the decoder):

$$\mathbf{Q} \in \mathbb{R}^{m \times d}$$

- A set of n **key** and n **value** vectors (come from the output of encoder):

$$\mathbf{K} \in \mathbb{R}^{n \times d}, \quad \mathbf{V} \in \mathbb{R}^{n \times d}.$$

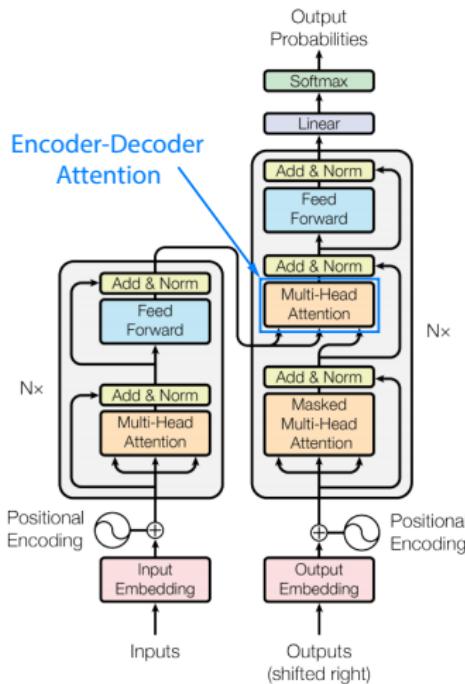
- Alignment scores and Compatibility Matrix **A**:

$$\mathbf{R} = \text{score}(\mathbf{Q}, \mathbf{K}) \in \mathbb{R}^{m \times n},$$

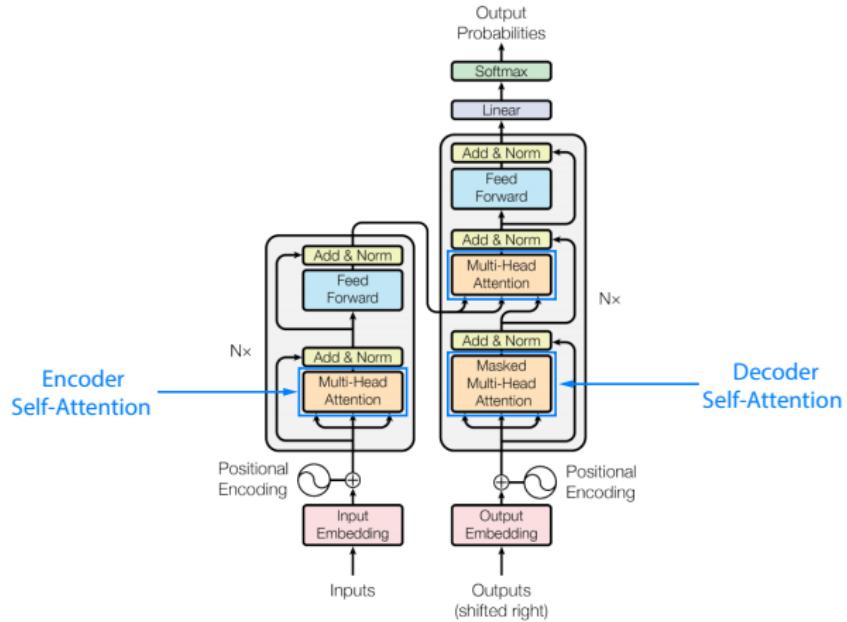
$$\mathbf{A} = \text{Softmax}(\mathbf{R}).$$

- Attention - weighted sum of value vectors:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \mathbf{A}\mathbf{V} \in \mathbb{R}^{m \times d}.$$



- Self-attention layer: **keys**, **values** and **queries** come from the **same place**.



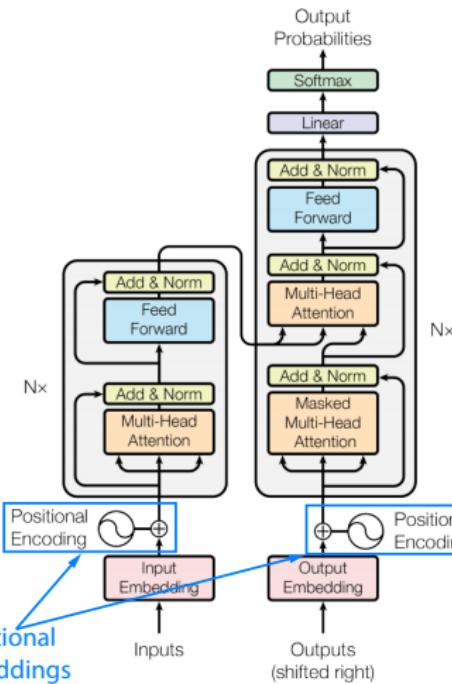
- Positional embeddings **PE**: vectors of the same dimension d as embeddings.
- Encode the positions of words within a sequence.
- Fixed positional embeddings:

$$PE_{(p,2i)} = \sin(p/10000^{2i/d}),$$

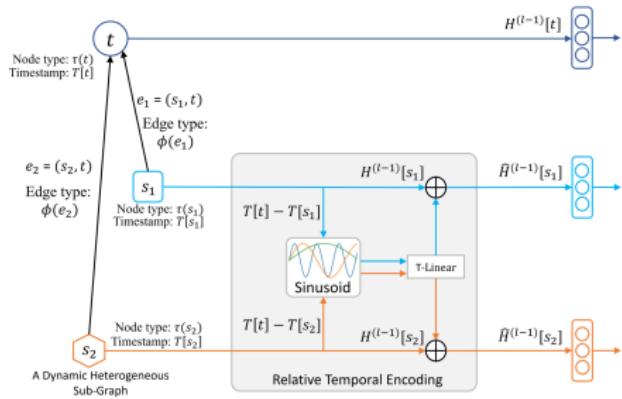
$$PE_{(p,2i+1)} = \cos(p/10000^{2i/d}).$$

- We don't need it in graphs: invariance to the order of nodes unless we are encoding some (time) dynamic.

Simply add the positional embedding to the word embedding.



- Handles graph dynamic



Heterogeneous Graph Transformer

$$\text{Base}(\Delta T(t, s), 2i) = \sin\left(\Delta T_{t,s}/10000^{\frac{2i}{d}}\right) \quad (6)$$

$$\text{Base}(\Delta T(t, s), 2i + 1) = \cos\left(\Delta T_{t,s}/10000^{\frac{2i+1}{d}}\right) \quad (7)$$

$$RTE(\Delta T(t, s)) = \text{T-Linear}\left(\text{Base}(\Delta T_{t,s})\right) \quad (8)$$

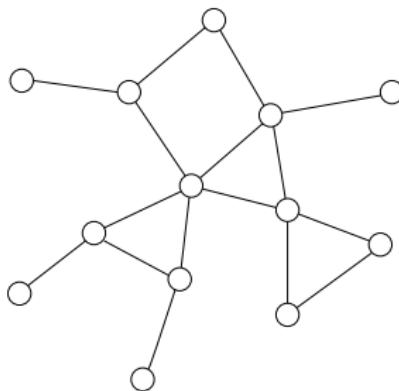
Finally, the temporal encoding relative to the target node t is added to the source node s ' representation as follows:

$$\tilde{H}^{(l-1)}[s] = H^{(l-1)}[s] + RTE(\Delta T(t, s)) \quad (9)$$

Graph Attention Networks

Graph Attention Networks
[arXiv:1710.10903 \[stat.ML\]](https://arxiv.org/abs/1710.10903)

- G - graph with N nodes and n features.



- Represent the graph in the tensor form:

$$\mathbf{G} \in \mathbb{R}^{N \times n}$$

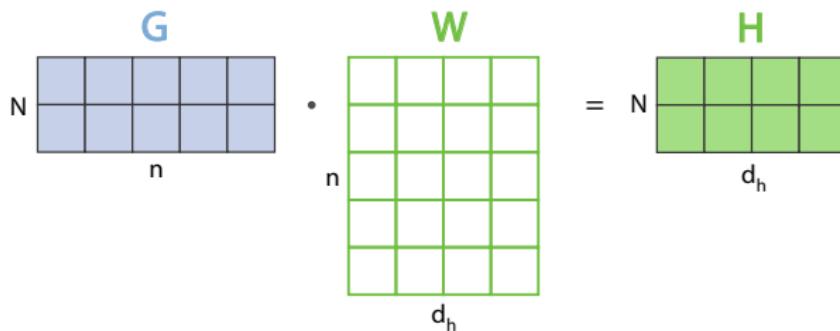
\mathbf{G}

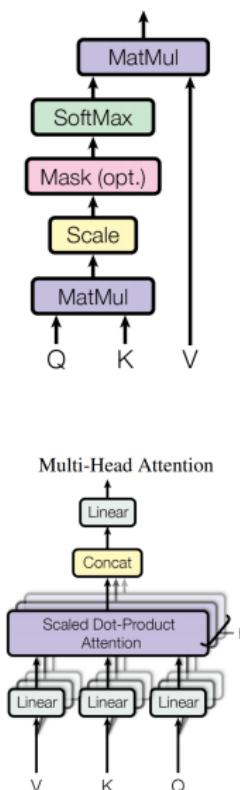
N	n
---	---

- Compute initial node embeddings through a linear projection $\mathbf{W} \in \mathbb{R}^{n \times d_h}$:

$$\mathbf{G} \in \mathbb{R}^{N \times n} \longrightarrow \mathbf{H} \in \mathbb{R}^{N \times d_h}$$

$$\mathbf{H} = \mathbf{G}\mathbf{W}$$





- Linearly project initial node embeddings $\mathbf{H} \in \mathbb{R}^{N \times d_h}$ (query, key, value):

$$\mathbf{Q} = \mathbf{H}\mathbf{W}^Q, \mathbf{K} = \mathbf{H}\mathbf{W}^K, \mathbf{V} = \mathbf{H}\mathbf{W}^V,$$

$$\mathbf{W}^Q, \mathbf{W}^K, \mathbf{W}^V \in \mathbb{R}^{d_h \times d}.$$

- Split into M heads and compute compatibility matrix $\mathbf{A} \in \mathbb{R}^{M \times N \times N}$ for graph nodes:

$$\mathbf{A} = \text{Softmax} \left(\frac{\mathbf{R}}{\sqrt{d/M}} \right) = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d/M}} \right).$$

- Compute attention messages for each head:

$$\mathbf{H}' = \mathbf{AV} \in \mathbb{R}^{M \times N \times d/M}$$

- Concatenate heads and project out with $\mathbf{W}^O \in \mathbb{R}^{d \times d}$:

$$MHA = \text{Concat}(\mathbf{H}'_1, \dots, \mathbf{H}'_M)\mathbf{W}^O \in \mathbb{R}^{N \times d}.$$

Multihead self-attention: step 1

Create query, key and value: $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$.

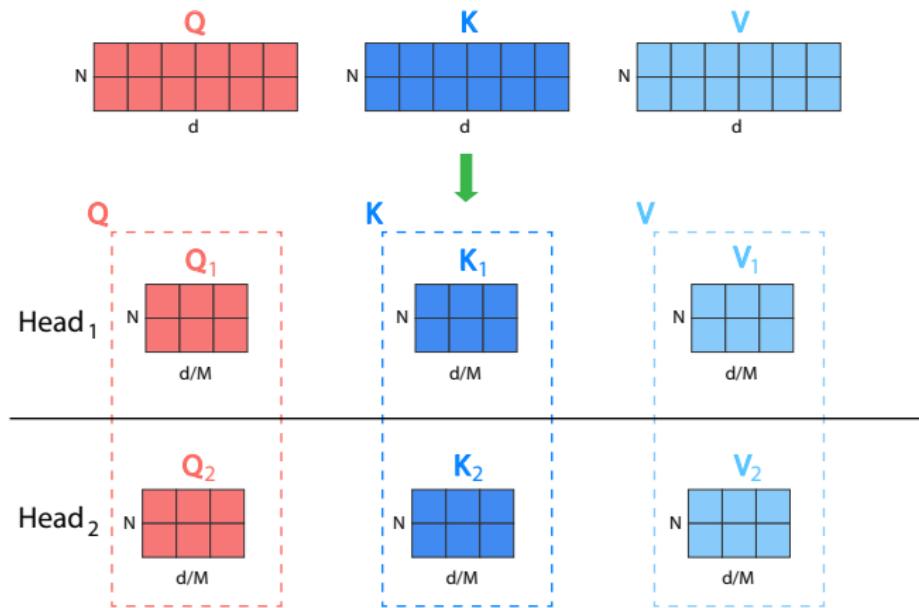
$$\begin{matrix} \text{H} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d_h & \end{matrix} \cdot \begin{matrix} \text{W}^Q \\ \text{---} \\ d_h & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix} = \begin{matrix} \text{Q} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix}$$

$$\begin{matrix} \text{H} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d_h & \end{matrix} \cdot \begin{matrix} \text{W}^K \\ \text{---} \\ d_h & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix} = \begin{matrix} \text{K} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix}$$

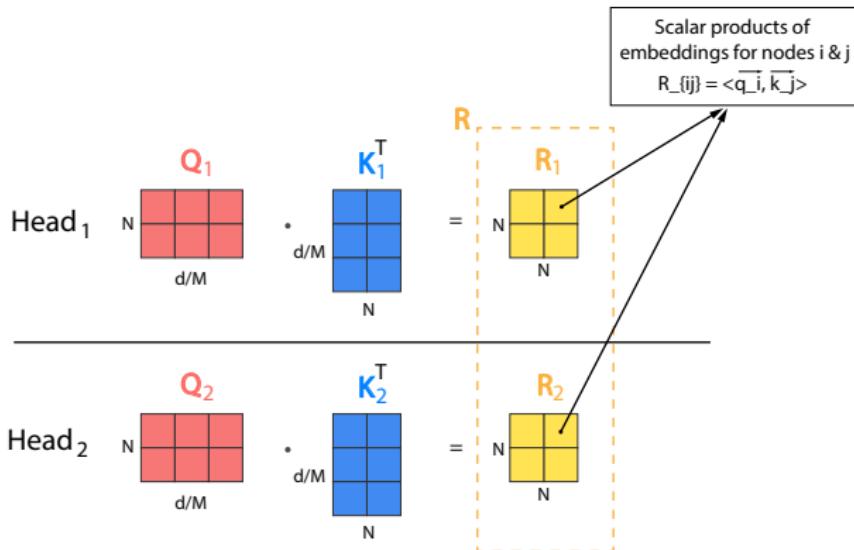
$$\begin{matrix} \text{H} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d_h & \end{matrix} \cdot \begin{matrix} \text{W}^V \\ \text{---} \\ d_h & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix} = \begin{matrix} \text{V} \\ \text{---} \\ N & \begin{matrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{matrix} \\ & d & \end{matrix}$$

Split computations into $M = 2$ heads.

This adds extra dimension: $[batch_size, N, d] \rightarrow [batch_size, M, N, d/M]$.



- Compute compatibility matrix for nodes: $\mathbf{A} \in \mathbb{R}^{M \times N \times N}$.
- Note how each head learns its own representations (*compare with convolutions*)

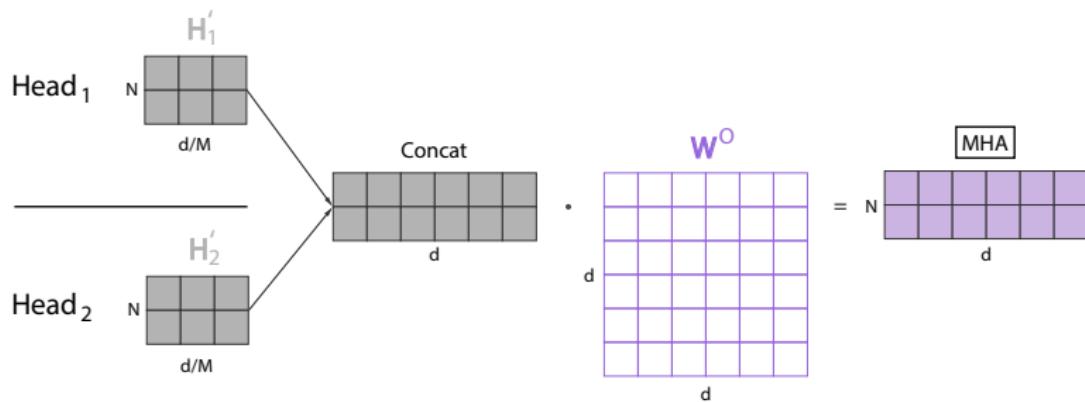


- Compute messages (attention outputs) for each head in parallel.
- Note how different heads update embeddings for each of N nodes using different subspaces of dimension d/M .

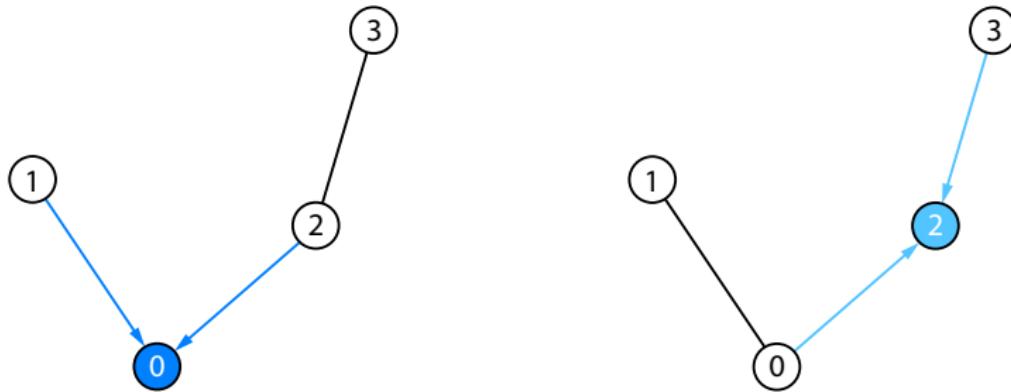
Head₁ Softmax $\left(\begin{matrix} \text{R}_1 \\ \hline N & \sqrt{d/M} \end{matrix} \right) \cdot \begin{matrix} \text{V}_1 \\ \hline N & d/M \end{matrix} = \begin{matrix} \text{H}'_1 \\ \hline N & d/M \end{matrix}$

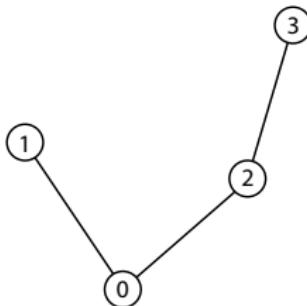
Head₂ Softmax $\left(\begin{matrix} \text{R}_2 \\ \hline N & \sqrt{d/M} \end{matrix} \right) \cdot \begin{matrix} \text{V}_2 \\ \hline N & d/M \end{matrix} = \begin{matrix} \text{H}'_2 \\ \hline N & d/M \end{matrix}$

- Concatenate outputs and project them to get the final values.



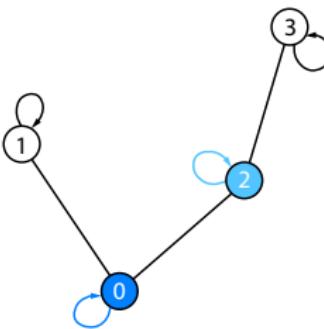
- How can we incorporate graph structure?
- We need to take into account only nodes in some neighbourhood (say, 1st order neighbours)





- Compute adjacency matrix:

$$\text{Adj} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix} \end{matrix}$$



- Add self-loops \iff add identity matrix:

$$\text{Adj} \longrightarrow \text{Adj} + I = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \left[\begin{matrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{matrix} \right] \end{matrix}$$

- Add mask:

$$\text{Mask} = -10^9 \cdot (1 - \text{Adj})$$

- Mask is applied before **Softmax**. Multiplied by (almost) negative infinity to zero out nodes after applying softmax.

$$\text{Adj} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \end{matrix}$$

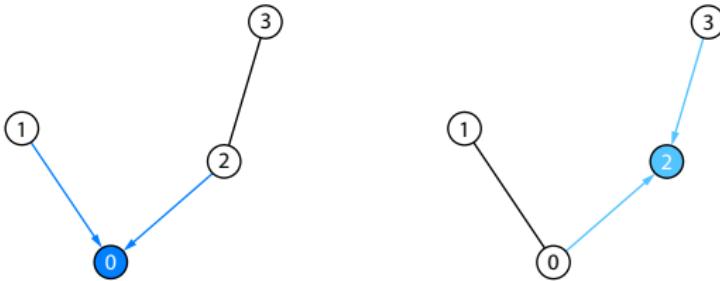
$$\text{Mask} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} 0 & 0 & 0 & -\text{np.inf} \\ 0 & 0 & -\text{np.inf} & -\text{np.inf} \\ 0 & -\text{np.inf} & 0 & 0 \\ -\text{np.inf} & -\text{np.inf} & 0 & 0 \end{bmatrix} \end{matrix}$$

- Suppose we have score matrix: $\mathbf{R} = \text{score}(\mathbf{Q}, \mathbf{K})$.
- Apply mask before computing attention messages:

$$\mathbf{R}' = \mathbf{R} + \text{Mask}$$

$$\mathbf{R}' = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} r_{00} & r_{01} & r_{02} & r_{03} \\ r_{10} & r_{11} & r_{12} & r_{13} \\ r_{20} & r_{21} & r_{22} & r_{23} \\ r_{30} & r_{31} & r_{32} & r_{33} \end{bmatrix} \end{matrix} + \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ 1 & \begin{bmatrix} 0 & 0 & 0 & -\text{np.inf} \\ 0 & 0 & -\text{np.inf} & -\text{np.inf} \\ 0 & -\text{np.inf} & 0 & 0 \\ -\text{np.inf} & -\text{np.inf} & 0 & 0 \end{bmatrix} \end{matrix}$$

$$\mathbf{A} = \text{Softmax}(\mathbf{R}') = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \begin{bmatrix} r_{00} & r_{01} & r_{02} & 0 \\ r_{10} & r_{11} & 0 & 0 \\ r_{20} & 0 & r_{22} & r_{23} \\ 0 & 0 & r_{32} & r_{33} \end{bmatrix} \end{matrix}$$



$$\mathbf{A} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \end{matrix} & \left[\begin{matrix} r_{00} & r_{01} & r_{02} & 0 \\ r_{10} & r_{11} & 0 & 0 \\ r_{20} & 0 & r_{22} & r_{23} \\ 0 & 0 & r_{32} & r_{33} \end{matrix} \right] \end{matrix}$$

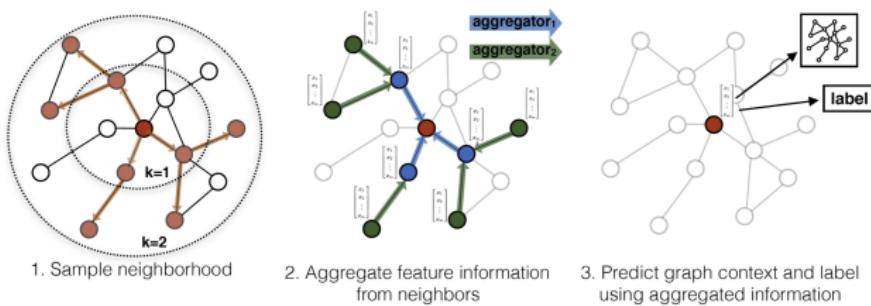
Only 1st order neighbours will send messages.

GraphSAGE

Inductive Representation Learning on Large Graphs

[arXiv:1706.02216 \[cs.SI\]](https://arxiv.org/abs/1706.02216)

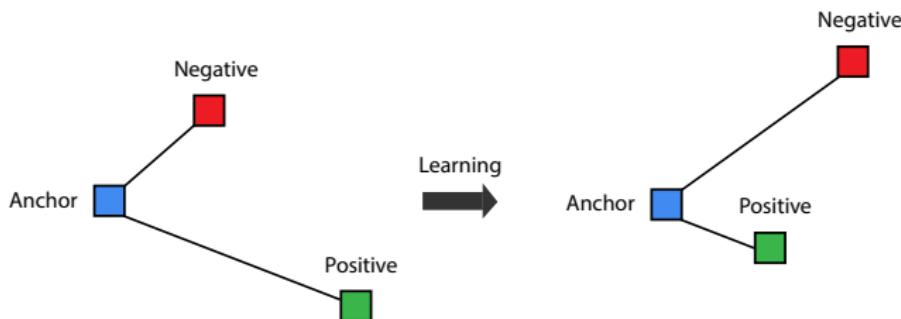
GraphSAGE: SAmples and aggreGatE



- Neighborhood sampling allows *mini-batch* training on large-scale graphs.
- Pytorch Geometric: [NeighborSampler](#) class.

Triplet Loss

- FaceNet: A Unified Embedding for Face Recognition and Clustering
[arXiv:1503.03832 \[cs.CV\]](https://arxiv.org/abs/1503.03832)
- Minimize the distance between an **anchor** and a **positive**.
- Maximize the distance between the **anchor** and a **negative**.

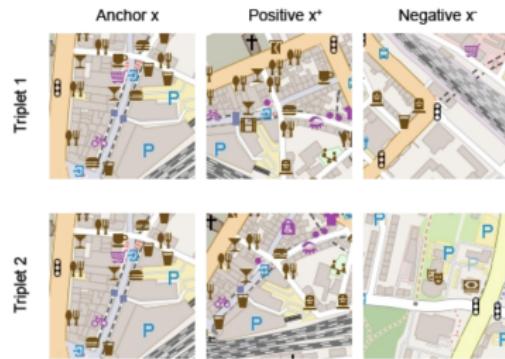


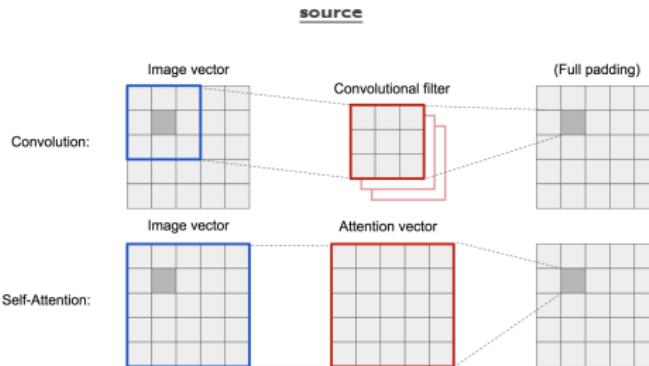
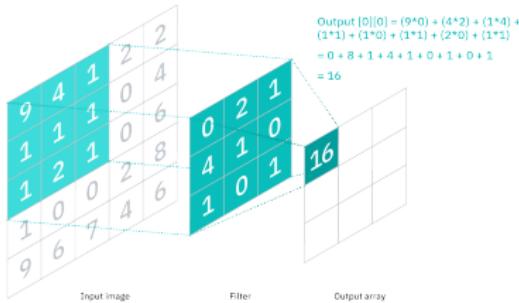
- Example: [loc2vec](#)

Location Embeddings

loc2vec

- N -channel tensor represents the area. Channels contain road network, points of interest, etc.
- Feature extractor: Convolutional neural network CNN

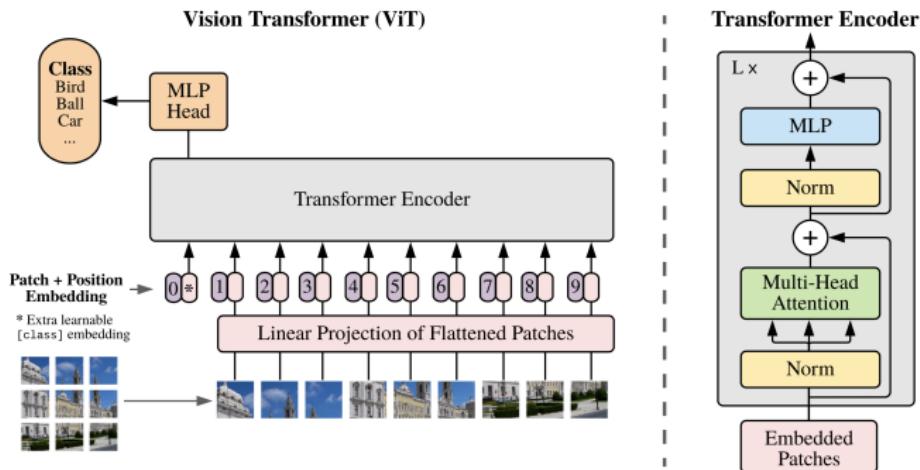


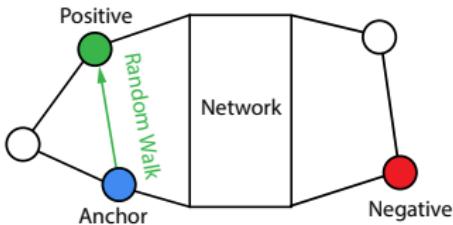


source

Vision Transformer

An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale
[arXiv:2010.11929 \[cs.CV\]](https://arxiv.org/abs/2010.11929)

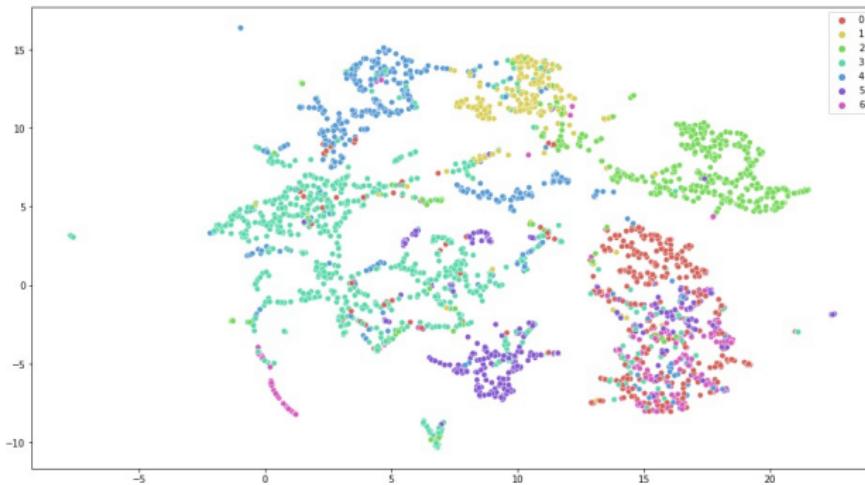




- Nearby nodes - similar representations
- Disparate nodes - distinct representations.

$$\mathcal{J}_G(\mathbf{z}_u) = -\log(\sigma(\mathbf{z}_u^T \mathbf{z}_v)) - Q \cdot \mathbb{E}_{v_n \sim P_n(v)} \log(\sigma(-\mathbf{z}_u^T \mathbf{z}_{v_n}))$$

2D visualization of CORA dataset unsupervised embeddings (GraphSAGE). Dimensional reduction was obtained by using [UMAP](#).



The End.