

DSC250: Advanced Data Mining

Topic Models / Text Embedding

Zhiting Hu

Lecture 9, Feb 4, 2025

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Recap: Expectation Maximization (EM)

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively:

E-step: $q(\mathbf{z} | \mathbf{x}) = p(\mathbf{z} | \mathbf{x}, \theta)$

We don't actually observe q , let's estimate it

M-step: $\max_{\theta} \mathbb{E}_{q(\mathbf{z} | \mathbf{x})} [\log p(\mathbf{x}, \mathbf{z} | \theta)]$

Let's "pretend" we also observe \mathbf{z} (its distribution)

Recap: Expectation Maximization (EM)

- Supervised MLE is easy:

$$\max_{\theta} \ell_c(\theta; \mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}, \mathbf{z} | \theta)$$

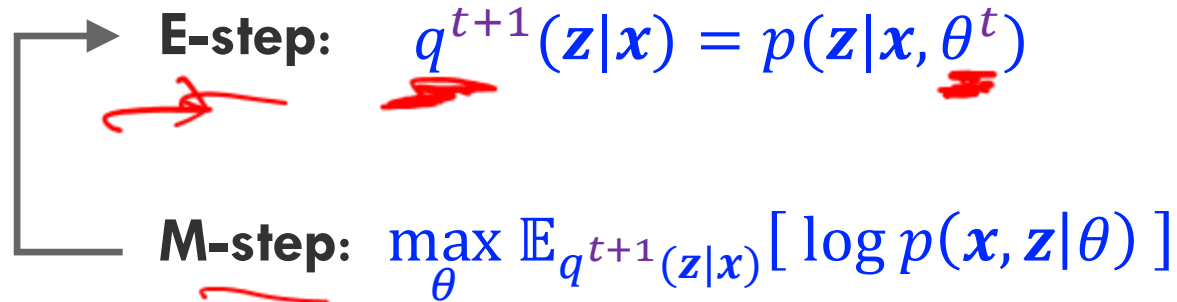
- Observe both \mathbf{x} and \mathbf{z}

- Unsupervised MLE is hard:

$$\max_{\theta} \ell(\theta; \mathbf{x}) = \log p(\mathbf{x} | \theta) = \log \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \theta)$$

- Observe only \mathbf{x}

- EM, intuitively:



We don't actually observe q , let's estimate it

Let's "pretend" we also observe \mathbf{z} (its distribution)

This is an iterative process

Recap: Expectation Maximization (EM)

- The EM algorithm is coordinate-decent on $F(q, \theta)$

- E-step: $q^{t+1} = \arg \min_q F(q, \theta^t) = p(\mathbf{z}|\mathbf{x}, \theta^t)$

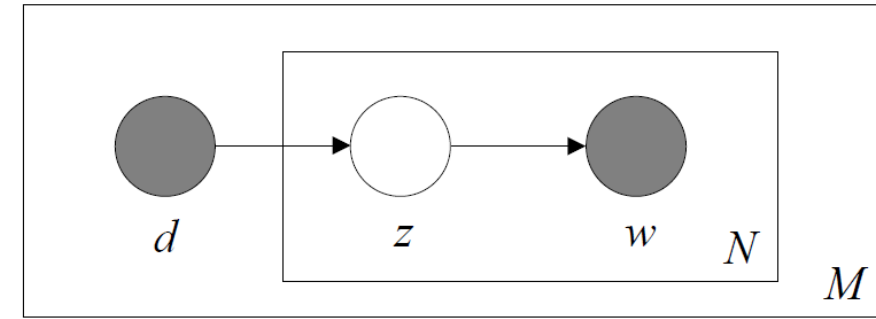
- the posterior distribution over the latent variables given the data and the current parameters

- M-step: $\theta^{t+1} = \arg \min_{\theta} F(q^{t+1}, \theta) = \operatorname{argmax}_{\theta} \sum_{\mathbf{z}} q^{t+1}(\mathbf{z}|\mathbf{x}) \log p(\mathbf{x}, \mathbf{z}|\theta)$

$\ell(\theta; \mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z}|\theta)}{q(\mathbf{z}|\mathbf{x})} \right] + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$

$= -F(q, \theta) + \text{KL}(q(\mathbf{z}|\mathbf{x}) || p(\mathbf{z}|\mathbf{x}, \theta))$

Recap: Learning pLSA with EM



- Likelihood function of a word w :

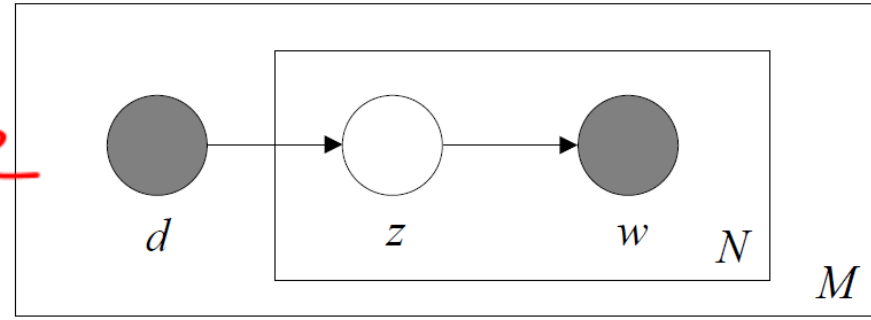
$$\begin{aligned} p(w|d, \theta, \beta) &= \sum_k p(w, z = k|d, \theta, \beta) \\ &= \sum_k p(w|z = k, d, \beta) p(z = k|d, \theta) = \sum_k \beta_{kw} \theta_{dk} \end{aligned}$$

- Learning by maximizing the log likelihood:

$$\max_{\beta, \theta} \log \sum_k \beta_{kw} \theta_{dk} \quad \checkmark$$

Recap: Learning pLSA with EM

Bayes rule



- E-step:

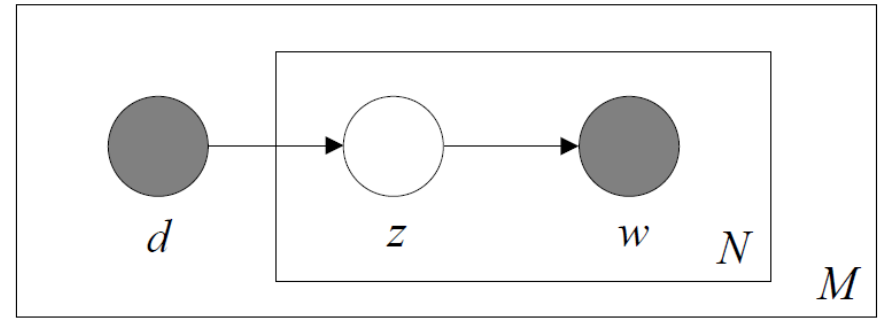
$$q(z|w,d) = P(z|w,d, \theta^t, \beta^t) = \frac{P(z,w,d, \theta^t, \beta^t)}{\sum_{z'} P(w,d, \theta^t, \beta^t)}$$

$$= \frac{\beta_{zw}^t \cdot \theta_{dz}^t}{\sum_{z'} \beta_{z'w}^t \cdot \theta_{dz'}^t}$$

- M-step:

$$\max_{\beta, \theta} E_{q(z|w,d)} \left[\log P(w,d, \theta^t, \beta^t) \right]$$

Recap: Learning pLSA with EM



- E-step:

$$p(z|w, d, \theta^t, \beta^t) = \frac{p(w|z, d, \beta^t)p(z|d, \theta^t)}{\sum_{z'} p(w|z', d, \beta^t)p(z'|d, \theta^t)} = \frac{\beta_{zw}^t \theta_{dz}^t}{\sum_{z'} \beta_{z'w}^t \theta_{dz'}^t}$$

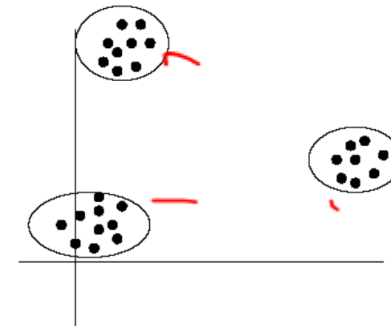
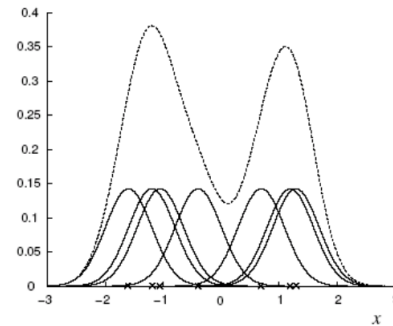
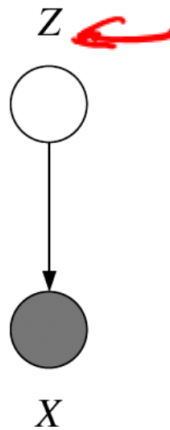
- M-step:

Another Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

$$p(x_n | \mu, \Sigma) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k)$$

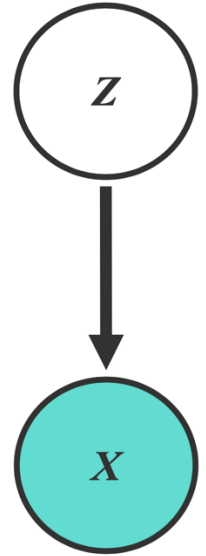
mixture proportion mixture component



$K=3$
200, 1, 2

- This model can be used for unsupervised clustering.
 - This model (fit by AutoClass) has been used to discover new kinds of stars in astronomical data, etc.

Example: Gaussian Mixture Models (GMMs)



- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

$$\approx \mathcal{N}(x_n | \mu_{z_n}, \Sigma_{z_n})$$

- The likelihood of a sample:

$$p(x | \pi, \mu, \Sigma) = \sum_z p(z | \pi) p(x | z, \mu, \Sigma)$$

$$\log p(x | \pi, \mu, \Sigma) = \log \sum_z \dots$$

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

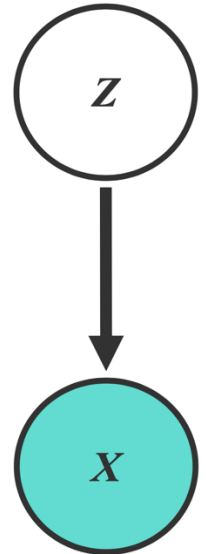
$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = \mathbf{1}, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

- The likelihood of a sample:

$$\begin{aligned} p(x_n | \mu, \Sigma) &= \sum_k p(z^k = \mathbf{1} | \pi) p(x, | z^k = \mathbf{1}, \mu, \Sigma) \\ &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k) \end{aligned}$$



Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components:

- Z is a latent class indicator vector:

$$p(z_n) = \text{multi}(z_n : \pi) = \prod_k (\pi_k)^{z_n^k}$$

- X is a conditional Gaussian variable with a class-specific mean/covariance

$$p(x_n | z_n^k = 1, \mu, \Sigma) = \frac{1}{(2\pi)^{m/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k)\right\}$$

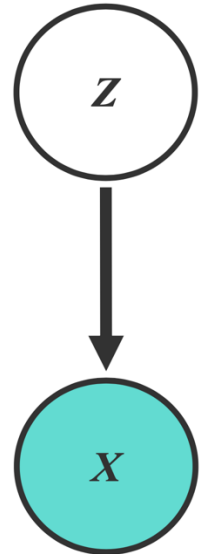
- The likelihood of a sample:

Parameters to be learned:

$$\begin{aligned}
 p(x_n | \mu, \Sigma) &= \sum_k p(z^k = 1 | \pi) p(x, | z^k = 1, \mu, \Sigma) \\
 &= \sum_{z_n} \prod_k \left((\pi_k)^{z_n^k} N(x_n : \mu_k, \Sigma_k)^{z_n^k} \right) = \sum_k \pi_k N(x, | \mu_k, \Sigma_k)
 \end{aligned}$$

mixture proportion

mixture component



Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$\begin{aligned}
 q(z|x, \mu^{(t)}, \Sigma^{(t)}, \pi^{(t)}) &= P(z|x, \mu^{(t)}, \Sigma^{(t)}, \pi^{(t)}) \\
 &= \frac{P(z, x, \mu^{(t)}, \Sigma^{(t)}, \pi^{(t)})}{\sum_{z'} P(z', x, \mu^{(t)}, \Sigma^{(t)}, \pi^{(t)})} \\
 &= \frac{\pi_z \cdot \mathcal{N}(x|z, \dots)}{\sum_{z'} \pi_{z'} \cdot \mathcal{N}(x|z', \dots)}
 \end{aligned}$$

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$p(z_n^k = 1 | x, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(x_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

Diagram illustrating the decomposition of the posterior probability:

- The numerator $\pi_k^{(t)} N(x_n | \mu_k^{(t)}, \Sigma_k^{(t)})$ is associated with the term $p(z_n^k = 1, x, \mu^{(t)}, \Sigma^{(t)})$.
- The denominator $\sum_i \pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})$ is associated with the term $p(x, \mu^{(t)}, \Sigma^{(t)})$.

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$p(z_n^k = 1 | x, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(x_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

$\nearrow p(z_n^k = 1, x, \mu^{(t)}, \Sigma^{(t)})$
 $\searrow p(x, \mu^{(t)}, \Sigma^{(t)})$

- M-step: the expected complete log likelihood

Handwritten: $\max_{\pi, \Sigma, \mu} E_{q(z | x, \mu^{(t)}, \Sigma^{(t)}, \pi^{(t)})} [\log P(z, x, \mu, \Sigma, \pi)]$

Example: Gaussian Mixture Models (GMMs)

- Consider a mixture of K Gaussian components
- E-step: computing the posterior of z_n given the current estimate of the parameters (i.e., π, μ, Σ)

$$p(z_n^k = 1 | x, \mu^{(t)}, \Sigma^{(t)}) = \frac{\pi_k^{(t)} N(x_n, | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_i \pi_i^{(t)} N(x_n, | \mu_i^{(t)}, \Sigma_i^{(t)})}$$

$\nearrow p(z_n^k = 1, x, \mu^{(t)}, \Sigma^{(t)})$
 $\searrow p(x, \mu^{(t)}, \Sigma^{(t)})$

- M-step: the expected complete log likelihood

$$\begin{aligned} \mathbb{E}_q [\ell_c(\theta; x, z)] &= \sum_n \mathbb{E}_q [\log p(z_n | \pi)] + \sum_n \mathbb{E}_q [\log p(x_n | z_n, \mu, \Sigma)] \\ &= \sum_n \sum_k \mathbb{E}_q [z_n^k] \log \pi_k - \frac{1}{2} \sum_n \sum_k \mathbb{E}_q [z_n^k] \left((x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \log |\Sigma_k| + C \right) \end{aligned}$$

Example: Gaussian Mixture Models (GMMs)

- M-step: computing the parameters given the current estimate of Z_n

$$\pi_k^* = \arg \max \langle l_c(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \frac{\partial}{\partial \pi_k} \langle l_c(\boldsymbol{\theta}) \rangle = 0, \forall k, \quad \text{s.t.} \quad \sum_k \pi_k = 1$$
$$\Rightarrow \quad \pi_k^* = \frac{\sum_n \langle Z_n^k \rangle_{q^{(t)}}}{N} = \frac{\sum_n \tau_n^{k(t)}}{N} = \frac{\langle n_k \rangle}{N}$$

$$\mu_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \mu_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} x_n}{\sum_n \tau_n^{k(t)}}$$

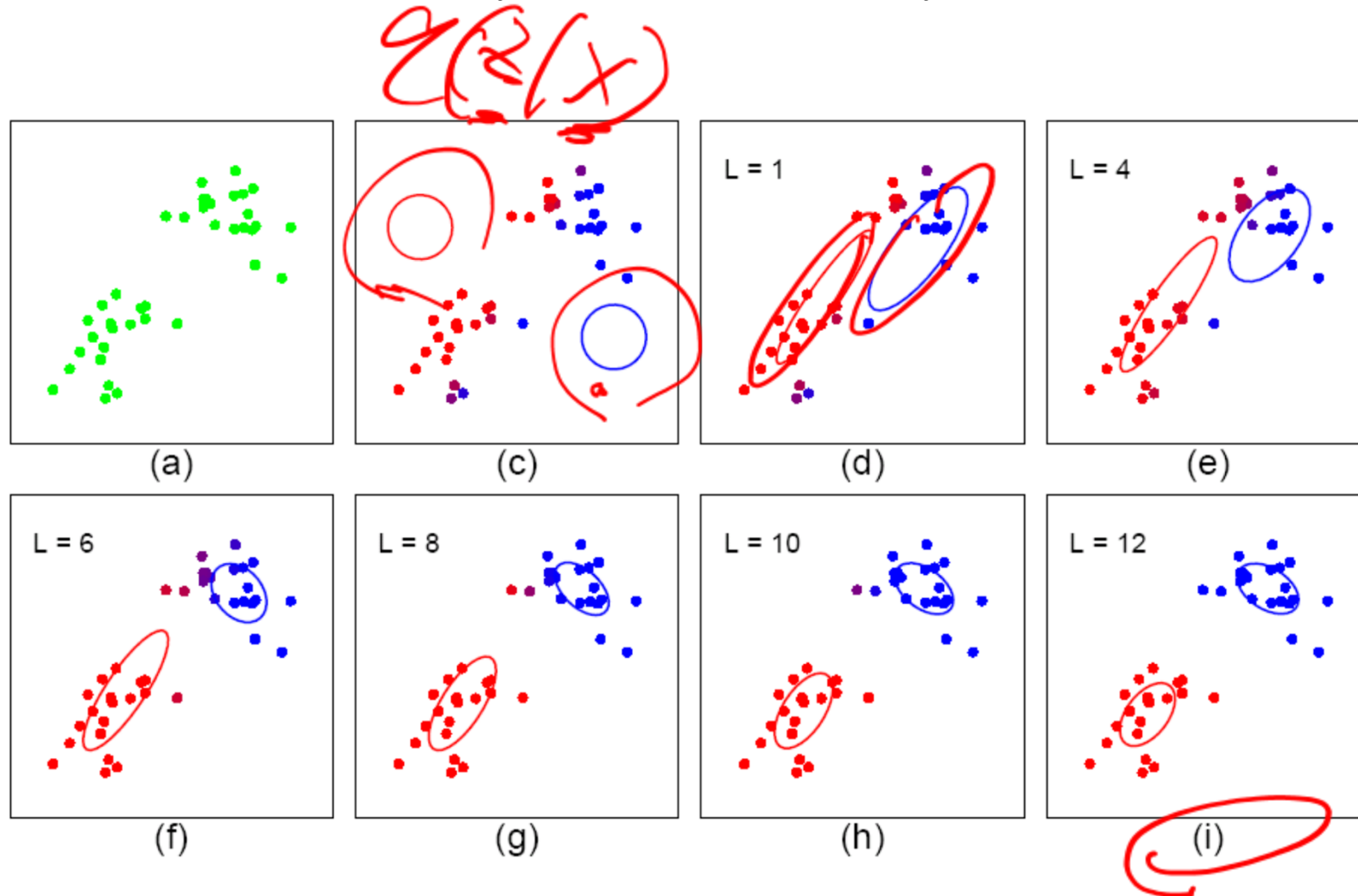
$$\Sigma_k^* = \arg \max \langle l(\boldsymbol{\theta}) \rangle, \quad \Rightarrow \quad \Sigma_k^{(t+1)} = \frac{\sum_n \tau_n^{k(t)} (x_n - \mu_k^{(t+1)})(x_n - \mu_k^{(t+1)})^T}{\sum_n \tau_n^{k(t)}}$$

Fact:

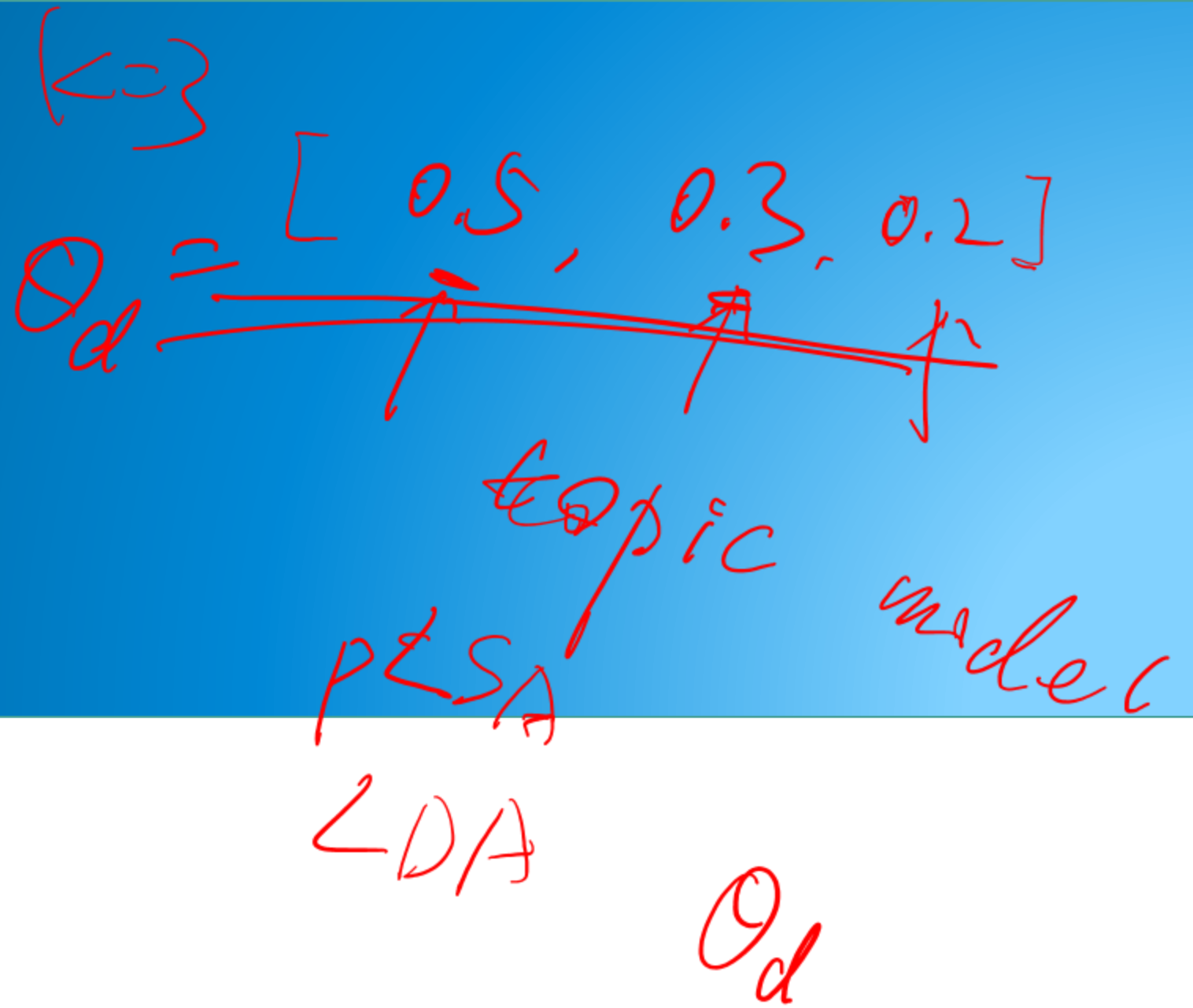
$$\frac{\partial \log |\mathbf{A}^{-1}|}{\partial \mathbf{A}^{-1}} = \mathbf{A}^T$$
$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{A}} = \mathbf{x} \mathbf{x}^T$$

Example: Gaussian Mixture Models (GMMs)

- Start: “guess” the centroid μ_k and covariance Σ_k of each of the K clusters
- Loop:



Text Embedding



Word Embedding

(distributed repr.) King - Man + Woman = Queen

- Conventional word embedding:
 - Word2vec, Glove
 - A pre-trained matrix, each row is an embedding vector of a word

	0	1	2	3	4	5	6	7	8	9
fox	-0.348680	-0.077720	0.177750	-0.094953	-0.452890	0.237790	0.209440	0.037886	0.035064	0.899010
ham	-0.773320	-0.282540	0.580760	0.841480	0.258540	0.585210	-0.021890	-0.463680	0.139070	0.658720
brown	-0.374120	-0.076264	0.109260	0.186620	0.029943	0.182700	-0.631980	0.133060	-0.128980	0.603430
beautiful	0.171200	0.534390	-0.348540	-0.097234	0.101800	-0.170860	0.295650	-0.041816	-0.516550	2.117200
jumps	-0.334840	0.215990	-0.350440	-0.260020	0.411070	0.154010	-0.386110	0.206380	0.386700	1.460500
eggs	-0.417810	-0.035192	-0.126150	-0.215930	-0.669740	0.513250	-0.797090	-0.068611	0.634660	1.256300
beans	-0.423290	-0.264500	0.200870	0.082187	0.066944	1.027600	-0.989140	-0.259950	0.145960	0.766450
sky	0.312550	-0.303080	0.019587	-0.354940	0.100180	-0.141530	-0.514270	0.886110	-0.530540	1.556600
bacon	-0.430730	-0.016025	0.484620	0.101390	-0.299200	0.761820	-0.353130	-0.325290	0.156730	0.873210
breakfast	0.073378	0.227670	0.208420	-0.456790	-0.078219	0.601960	-0.024494	-0.467980	0.054627	2.283700
toast	0.130740	-0.193730	0.253270	0.090102	-0.272580	-0.030571	0.096945	-0.115060	0.484000	0.848380
today	-0.156570	0.594890	-0.031445	-0.077586	0.278630	-0.509210	-0.066350	-0.081890	-0.047986	2.803600
blue	0.129450	0.036518	0.032298	-0.060034	0.399840	-0.103020	-0.507880	0.076630	-0.422920	0.815730
green	-0.072368	0.233200	0.137260	-0.156630	0.248440	0.349870	-0.241700	-0.091426	-0.530150	1.341300
kings	0.259230	-0.854690	0.360010	-0.642000	0.568530	-0.321420	0.173250	0.133030	-0.089720	1.528600
dog	-0.057120	0.052685	0.003026	-0.048517	0.007043	0.041856	-0.024704	-0.039783	0.009614	0.308416
sausages	-0.174290	-0.064869	-0.046976	0.287420	-0.128150	0.647630	0.056315	-0.240440	-0.025094	0.502220
lazy	-0.353320	-0.299710	-0.176230	-0.321940	-0.385640	0.586110	0.411160	-0.418680	0.073093	1.486500
love	0.139490	0.534530	-0.252470	-0.125650	0.048748	0.152440	0.199060	-0.065970	0.128830	2.055900
quick	-0.445630	0.191510	-0.249210	0.465900	0.161950	0.212780	-0.046480	0.021170	0.417660	1.686900

CDA
2003
2013
2023
low-pn base pape

[Courtesy: Vaswani, et al., 2017]

20 rows x 300 columns

Word Embedding

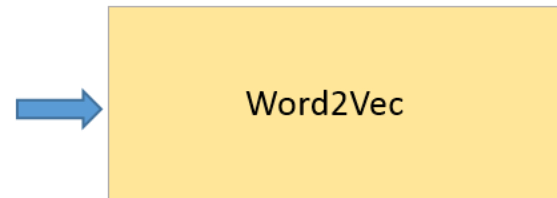
- Conventional word embedding:
 - Word2vec, Glove
 - A pre-trained **matrix**, each row is an embedding vector of a word

	0	1	2	3	4	5	6	7	8	9
fox	-0.348680	-0.077720	0.177750	-0.094953	-0.452890	0.237790	0.209440	0.037886	0.035064	0.899010
ham	-0.773320	-0.282540	0.580760	0.841480	0.258540	0.585210	-0.021890	-0.463680	0.139070	0.658720
brown	-0.374120	-0.076264	0.109260	0.186620	0.029943	0.182700	-0.631980	0.133060	-0.128980	0.603430
beautiful	0.171200	0.534390	-0.348540	-0.097234	0.101800	-0.170860	0.295650	-0.041816	-0.516550	2.117200
jumps	-0.334840	0.215990	-0.350440	-0.260020	0.411070	0.154010	-0.386110	0.206380	0.386700	1.460500
eggs	-0.417810	-0.035192	-0.126150	-0.215930	-0.669740	0.513250	-0.797090	-0.068611	0.634660	1.256300
beans	-0.423290	-0.264500	0.200870	0.082187	0.066944	1.027600	-0.989140	-0.259950	0.145960	0.766450
sky	0.312550	-0.303080	0.019587	-0.354940	0.100180	-0.141530	-0.514270	0.886110	-0.530540	1.556600
bacon	-0.430730	-0.016025	0.484620	0.101390	-0.299200	0.761820	-0.353130	-0.325290	0.156730	0.873210
breakfast	0.073378	0.227670	0.208420	-0.456790	-0.078219	0.601960	-0.024494	-0.467980	0.054627	2.283700

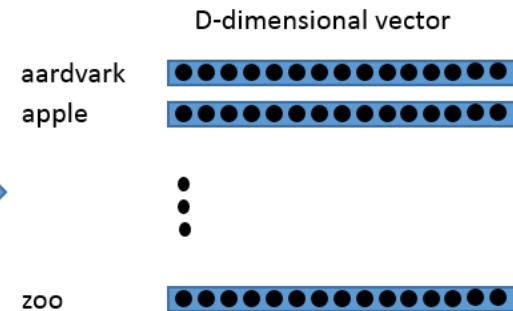
English Wikipedia Corpus

The Annual Reminder continued through July 4, 1969. This final Annual Reminder took place less than a week after the June 28 Stonewall riots, in which the patrons of the Stonewall Inn, a gay bar in Greenwich Village, fought against police who raided the bar. Rodwell received several telephone calls threatening him and the other New York participants, but he was able to arrange for police protection for the chartered bus all the way to Philadelphia. About 45 people participated, including the deputy mayor of Philadelphia and his wife. The dress code was still in effect at the Reminder, but two women from the New York contingent broke from the single-file picket line and held hands. When Kameny tried to break them apart, Rodwell furiously denounced him to onlooking members of the press.

Following the 1969 Annual Reminder, there was a sense, particularly among the younger and more radical participants, that the time for silent picketing had passed. Dissent and dissatisfaction had begun to take new and more emphatic forms in society. "The conference passed a resolution drafted by Rodwell, his partner Fred Sargeant, Broidy and Linda Rhodes to move the demonstration from July 4 in Philadelphia to the last weekend in June in New York City, as well as proposing to "other organizations throughout the country... suggesting that they hold parallel demonstrations on that day" to commemorate the Stonewall riot.



Embedding Matrix



3945	-0.115060	0.484000	0.848380
3350	-0.081890	-0.047986	2.803600
7880	0.076630	-0.422920	0.815730
1700	-0.091426	-0.530150	1.341300
3250	0.133030	-0.089720	1.528600
1704	-0.039783	0.009614	0.308416
3315	-0.240440	-0.025094	0.502220
1160	-0.418680	0.073093	1.486500
3060	-0.065970	0.128830	2.055900
3480	0.021170	0.417660	1.686900

Word Embedding

- **Problem:** word embeddings are applied in a context free manner

open a bank account

on the river bank

[0.3, 0.2, -0.8, ...]



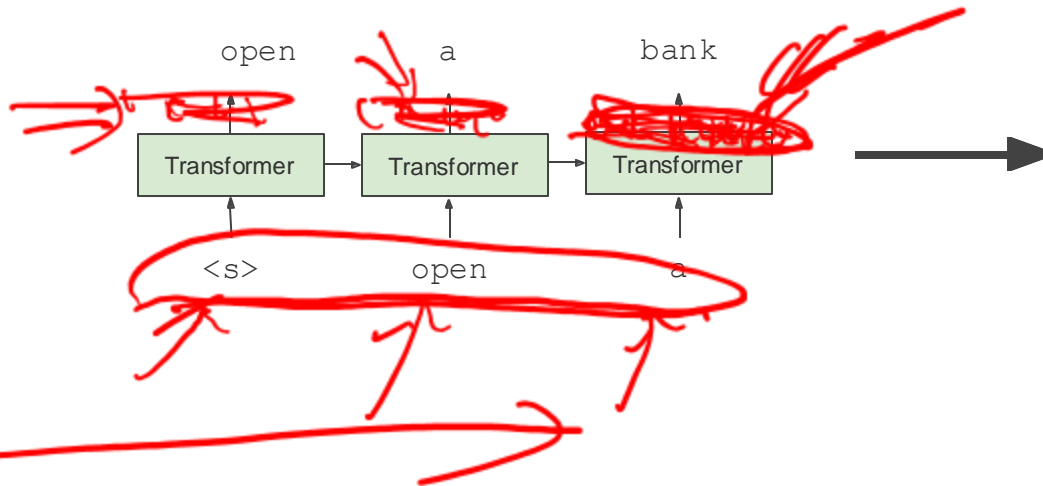
Contextual Representations

GPT-1

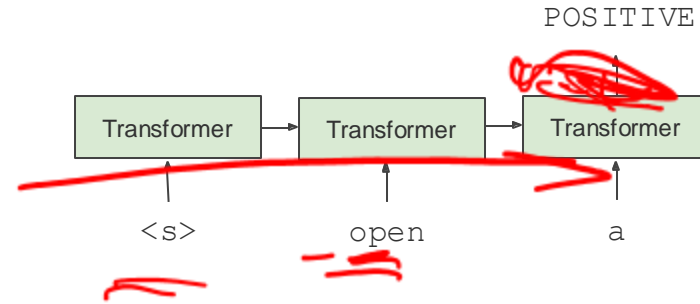
→ GPT-v1

- *Improving Language Understanding by Generative Pre-Training*, OpenAI, 2018

Train Deep (12-layer) Transformer LM

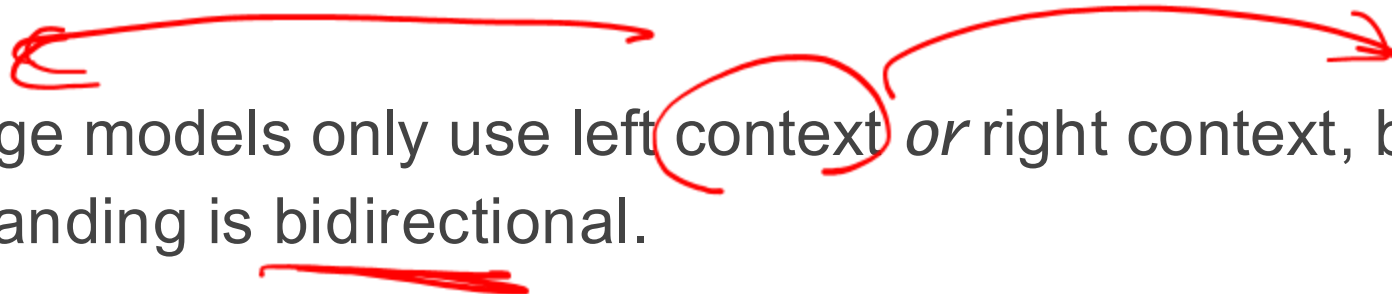


Fine-tune on Classification Task



Problem with Previous Methods

- **Problem:** Language models only use left context *or* right context, but language understanding is bidirectional.

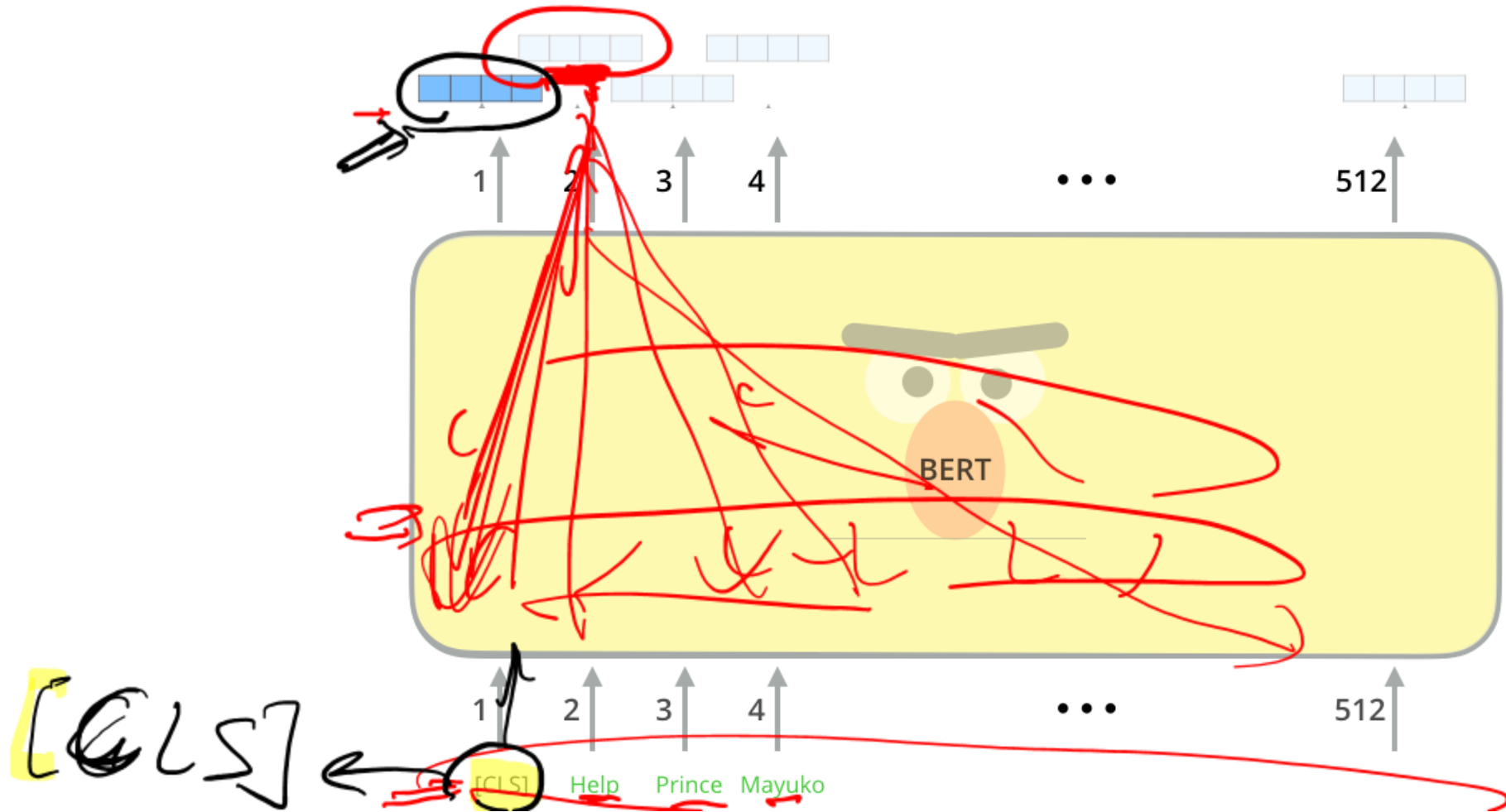


2018

Pretraining

BERT

- BERT: A bidirectional model to extract contextual word embedding



BERT: Pre-training Procedure

- Dataset:
 - Wikipedia (2.5B words) + a collection of free ebooks (800M words)

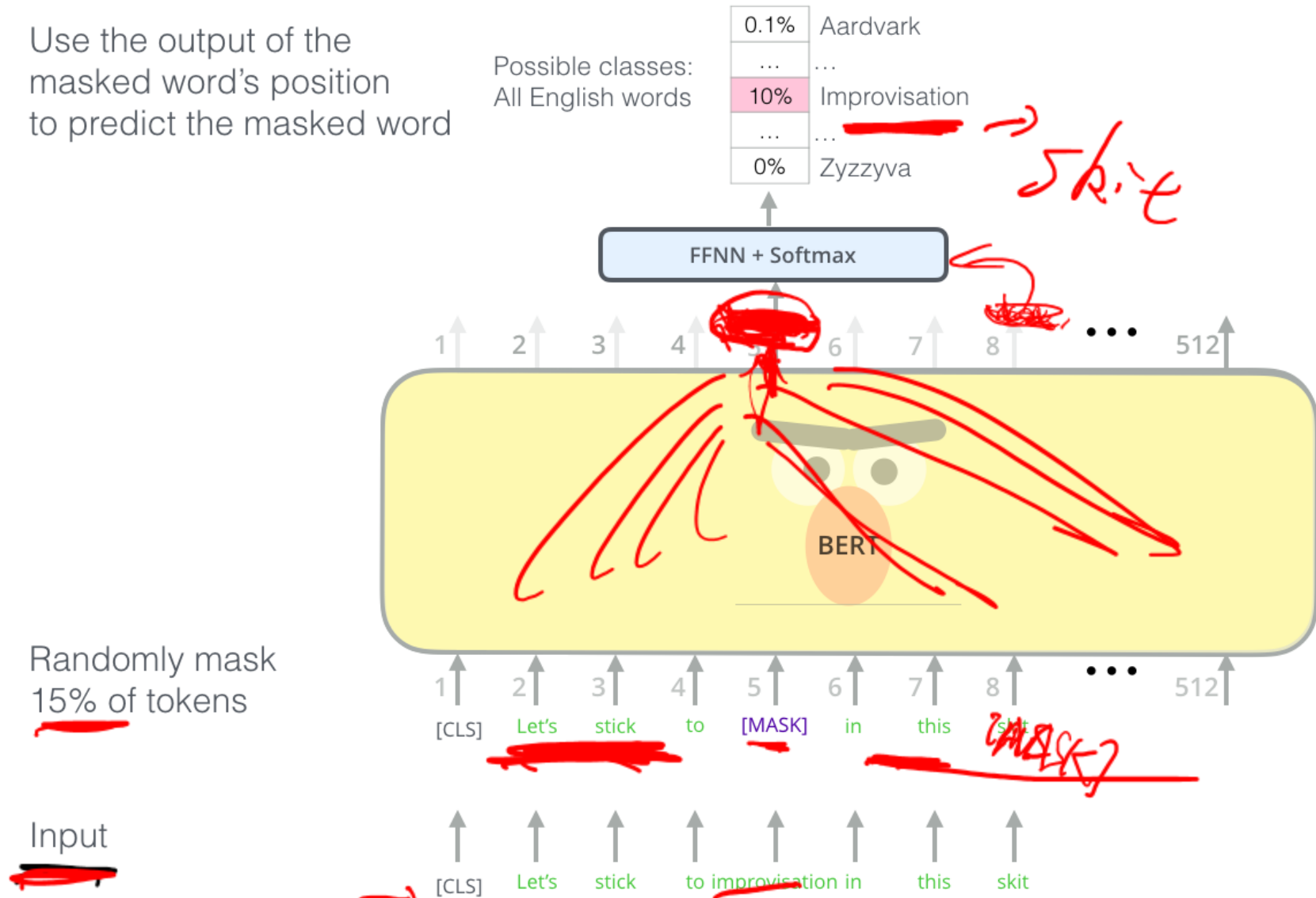
BERT: Pre-training Procedure

- Dataset:
 - Wikipedia (2.5B words) + a collection of free ebooks (800M words)
- Training procedure
 - **masked language model** (masked LM)
 - Masks some percent of words from the input and has to reconstruct those words from context

BERT: Pre-training Procedure

- Masked LM

Use the output of the masked word's position to predict the masked word



BERT: Pre-training Procedure

- Masked LM
- 15% masking:
 - Too little masking: Too expensive to train (few supervision signals per example)
 - Too much masking: Not enough context
- Problem: Mask token never seen at fine-tuning
- Solution: don't replace with [MASK] 100% of the time. Instead:
- 80% of the time, replace with [MASK]
 - went to the store → went to the [MASK]
- 10% of the time, replace random word
 - went to the store → went to the running
- 10% of the time, keep same
 - went to the store → went to the store



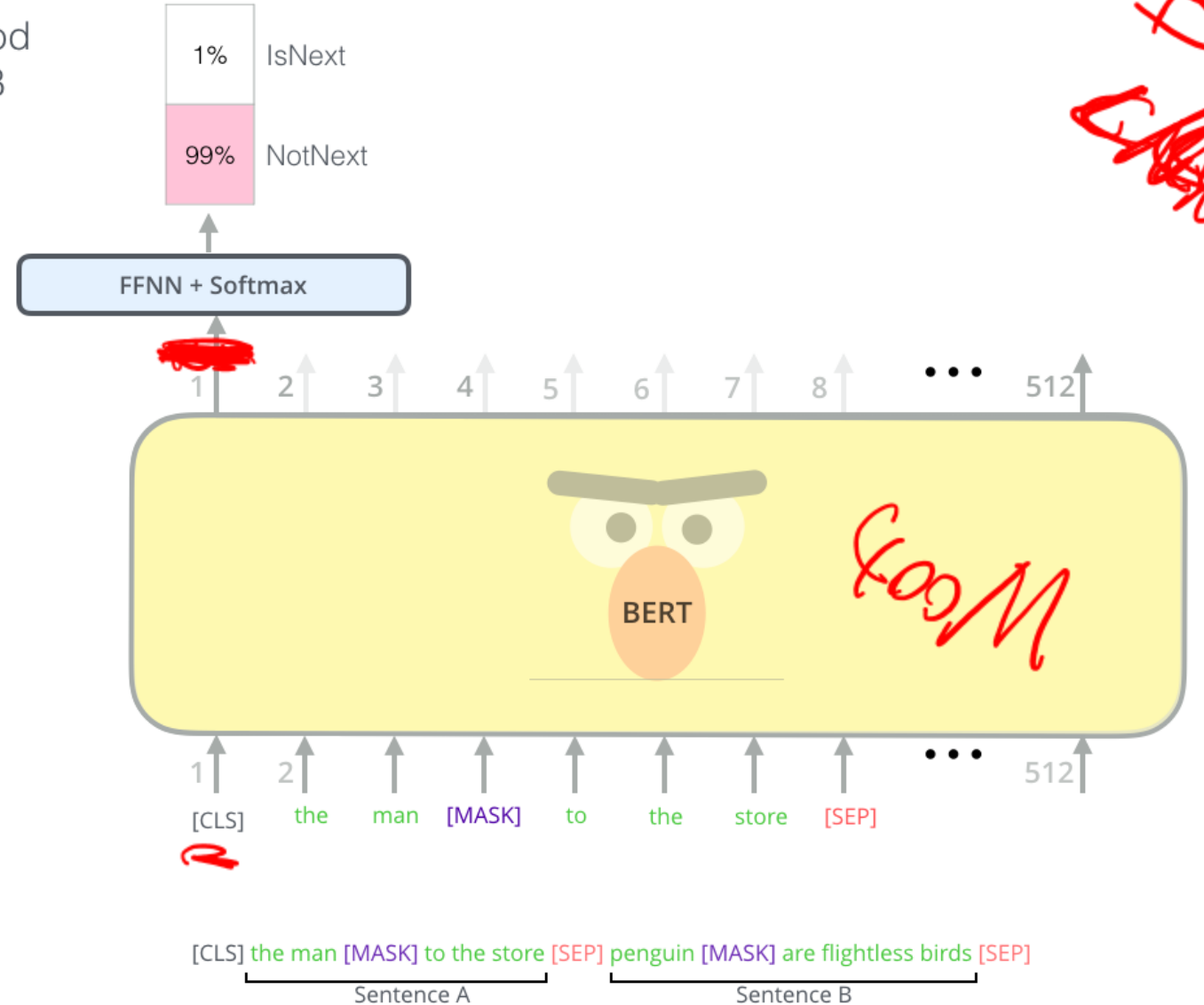
BERT: Pre-training Procedure

- Dataset:
 - Wikipedia (2.5B words) + a collection of free ebooks (800M words)
- Training procedure
 - **masked language model** (masked LM)
 - Masks some percent of words from the input and has to reconstruct those words from context
 - **Two-sentence task**
 - To understand relationships between sentences
 - Concatenate two sentences A and B and predict whether B actually comes after A in the original text

BERT: Pre-training Procedure

- Two sentence task

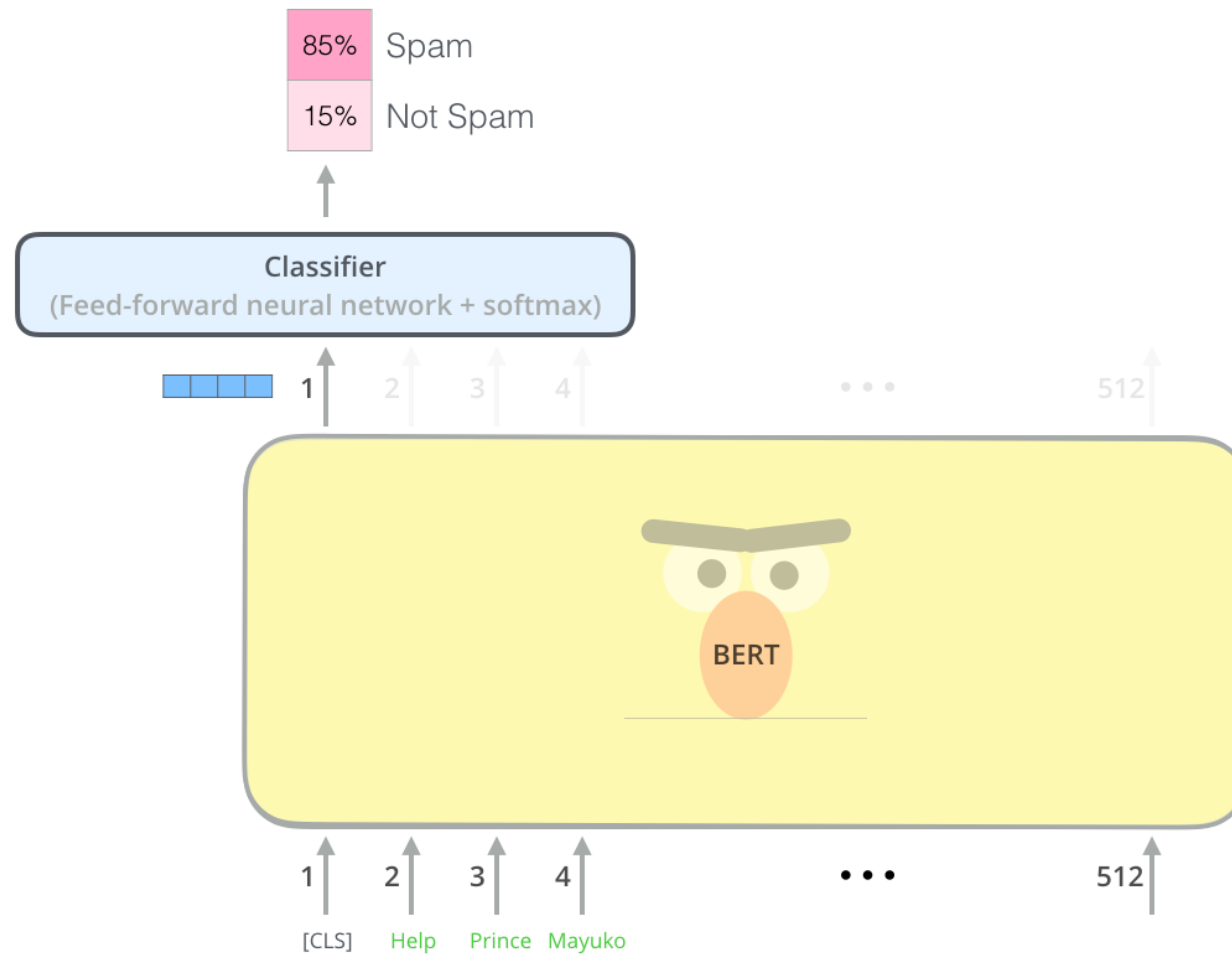
Predict likelihood that sentence B belongs after sentence A



[CLS] A
[SEP] B

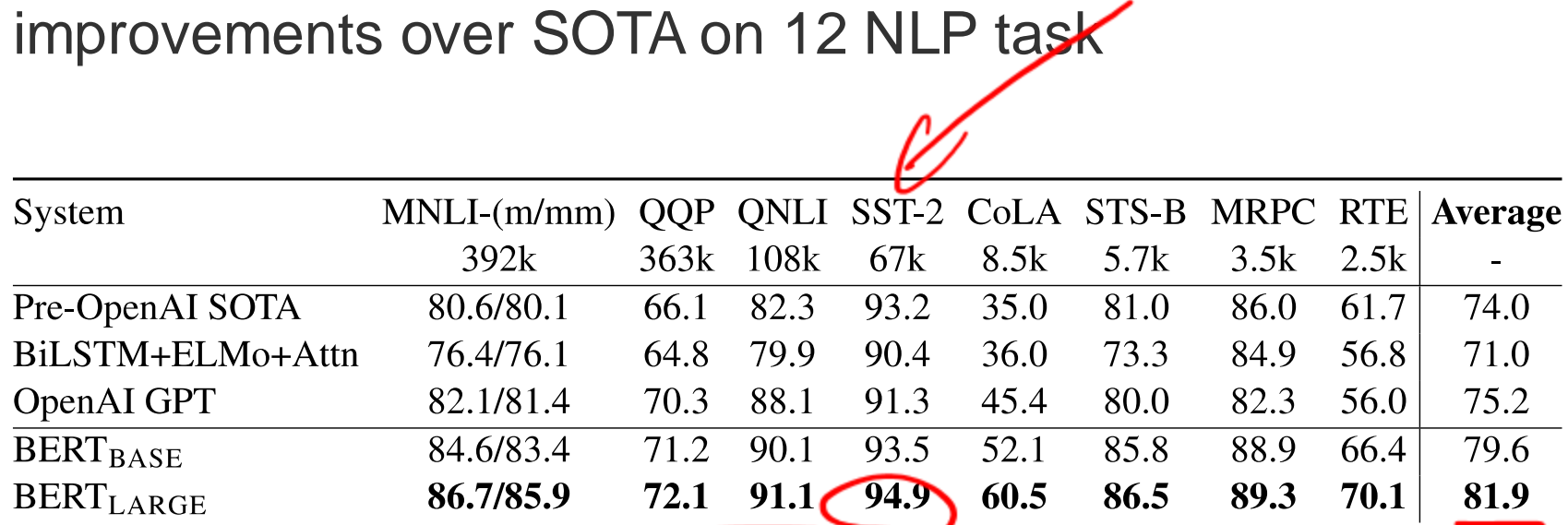
BERT: Downstream Fine-tuning

- Use BERT for sentence classification



BERT Results

- Huge improvements over SOTA on 12 NLP task



System	MNLI-(m/mm) 392k	QQP 363k	QNLI 108k	SST-2 67k	CoLA 8.5k	STS-B 5.7k	MRPC 3.5k	RTE 2.5k	Average
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERT _{LARGE}	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The “Average” column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT_{BASE} = (L=12, H=768, A=12); BERT_{LARGE} = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from <https://gluebenchmark.com/leaderboard> and <https://blog.openai.com/language-unsupervised/>.

Similar idea for image embedding: masked autoencoder (MAE)

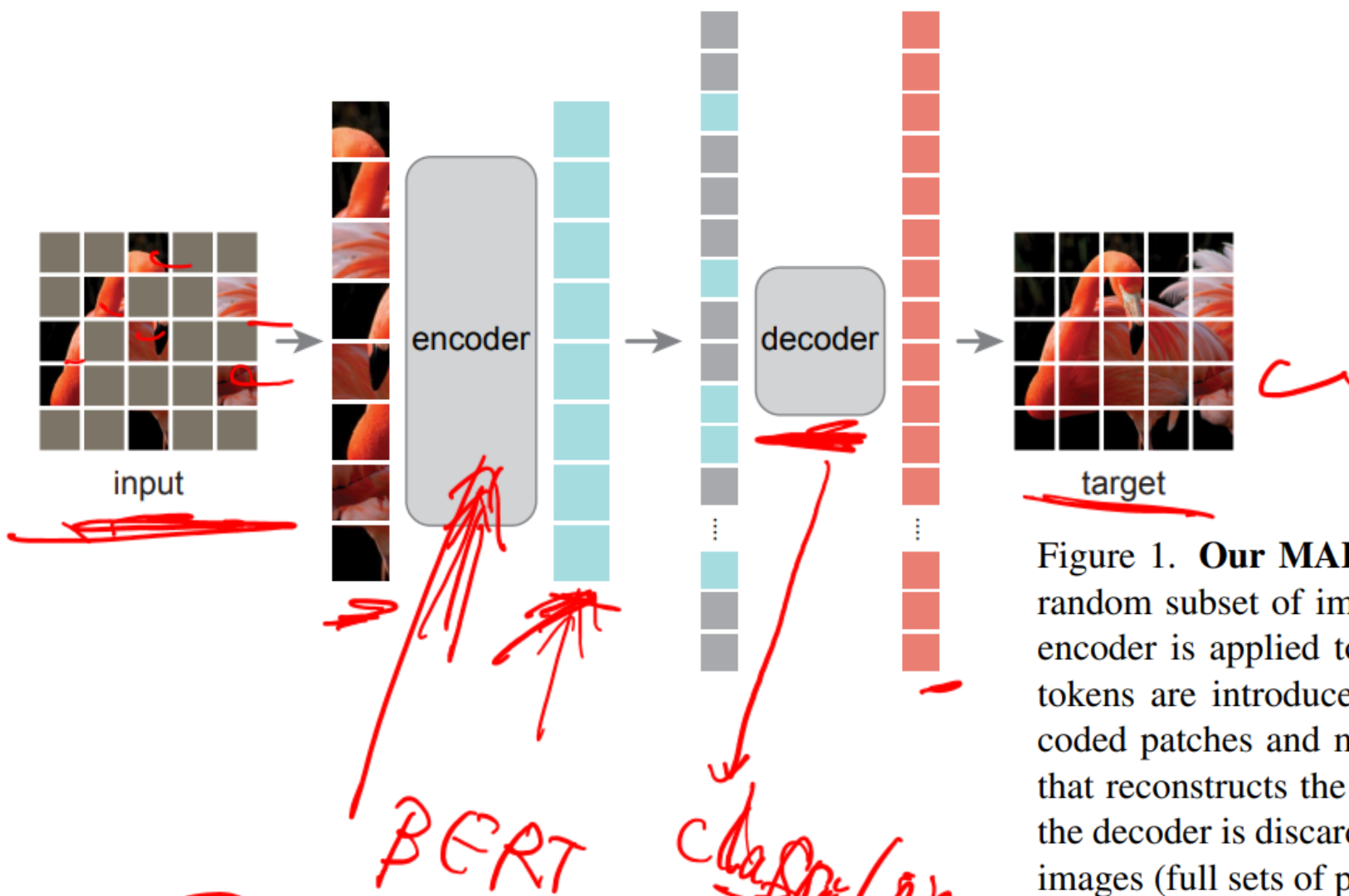
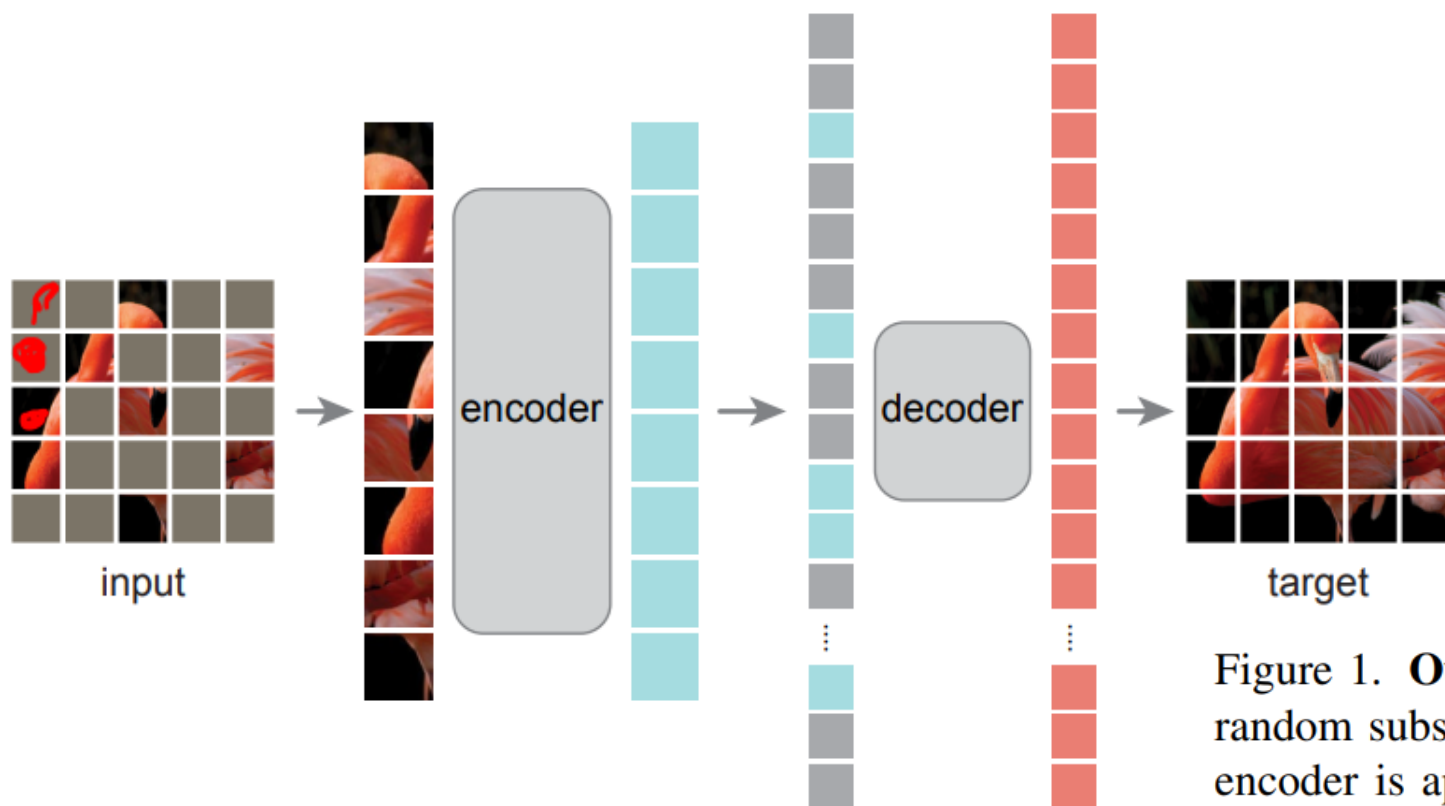


Figure 1. **Our MAE architecture.** During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of ~~visible patches~~. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

Similar idea for image embedding: masked autoencoder (MAE)

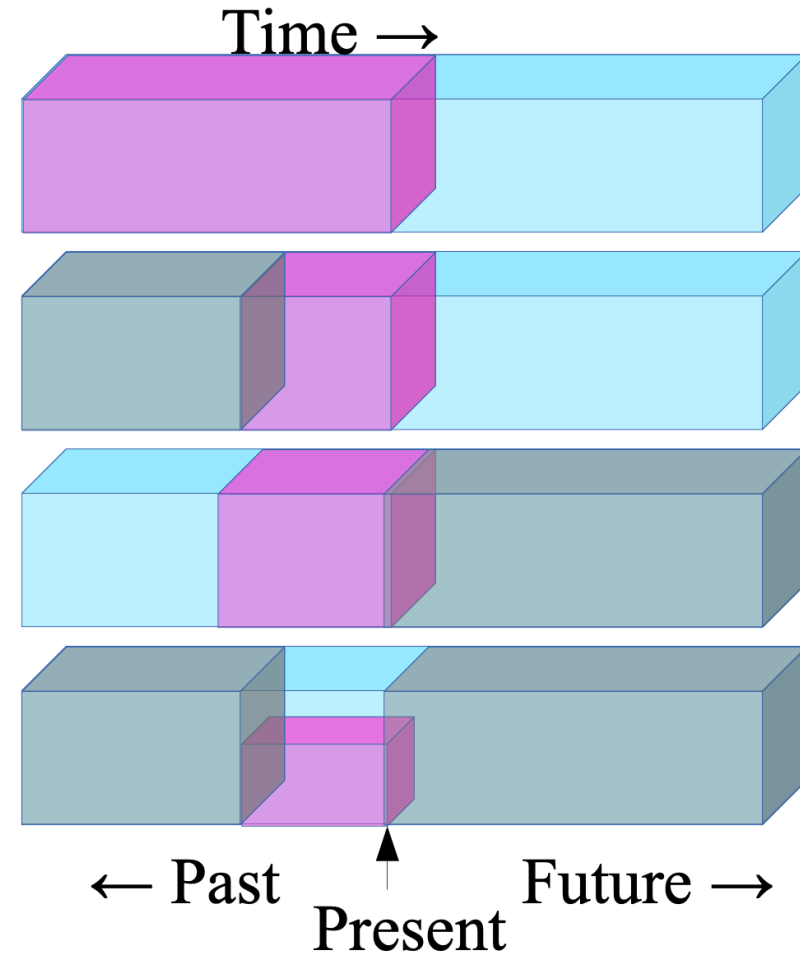


Question: Why is this (75%) much larger than the mask rate in BERT (15%)?

Figure 1. **Our MAE architecture.** During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

More general idea: Self-Supervised Learning

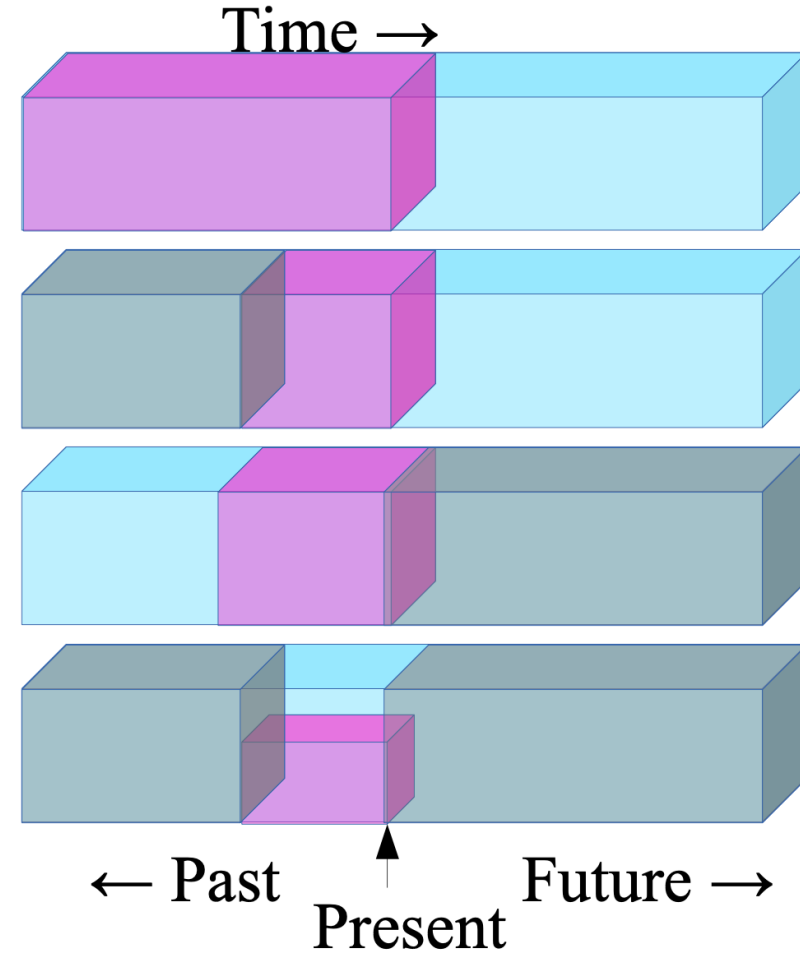
- ▶ Predict any part of the **input** from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.



2d/8
Self-supervised

More general idea: Self-Supervised Learning

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the ~~occluded~~ from the visible
- ▶ Pretend there is a part of the input you don't know and predict that.



More general idea: Self-Supervised Learning: Motivation (I)

- ▶ Our brains do this all the time
- ▶ Filling in the visual field at the retinal blind spot
- ▶ Filling in occluded images, missing segments in speech
- ▶ Predicting the state of the world from partial (textual) descriptions
- ▶ Predicting the consequences of our actions
- ▶ Predicting the sequence of actions leading to a result
- ▶ Predicting any part of the past, present or future percepts from whatever information is available.



More general idea: Self-Supervised Learning: Motivation (I)

- Successfully learning to predict everything from everything else would result in **the accumulation of lots of background knowledge about how the world works**
- The model is forced to learn what we really care about, e.g. a semantic representation, in order to solve the prediction problem

[Courtesy: Lecun “Self-supervised Learning”]

[Courtesy: Zisserman “Self-supervised Learning”]

More general idea: Self-Supervised Learning: Motivation (II)

- The machine predicts any part of its input from any observed part
 - A lot of supervision signals in each data instance
- Untapped/availability of vast numbers of unlabeled text/images/videos..
 - Facebook: one billion images uploaded per day
 - 300 hours of video are uploaded to YouTube every minute

Costs a million of each video

Questions?