
Foundations for Natural Language Processing

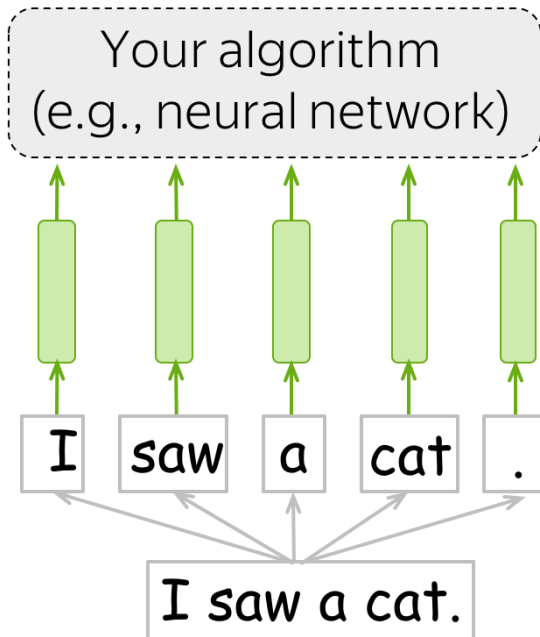
Neural Embeddings

Ivan Titov

(with graphics/materials from Elena Voita)



Neural models and word embeddings



Any algorithm for solving a task

Word representation - vector
(input for your model/algorithm)

Sequence of tokens

Text (your input)

Neural models and word embeddings

Token index in
the vocabulary

39 1592 10 2548 5



I

saw

a

cat

.

Embedding
dimension



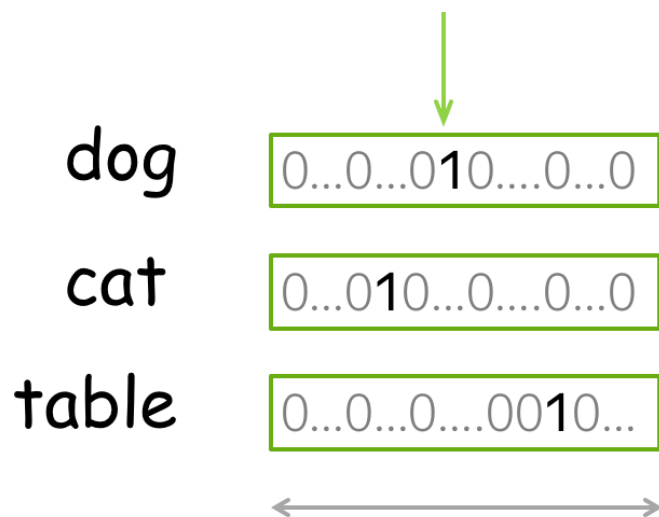
39



Vocabulary
size

One-hot vectors as word representations

One is 1, the rest are 0

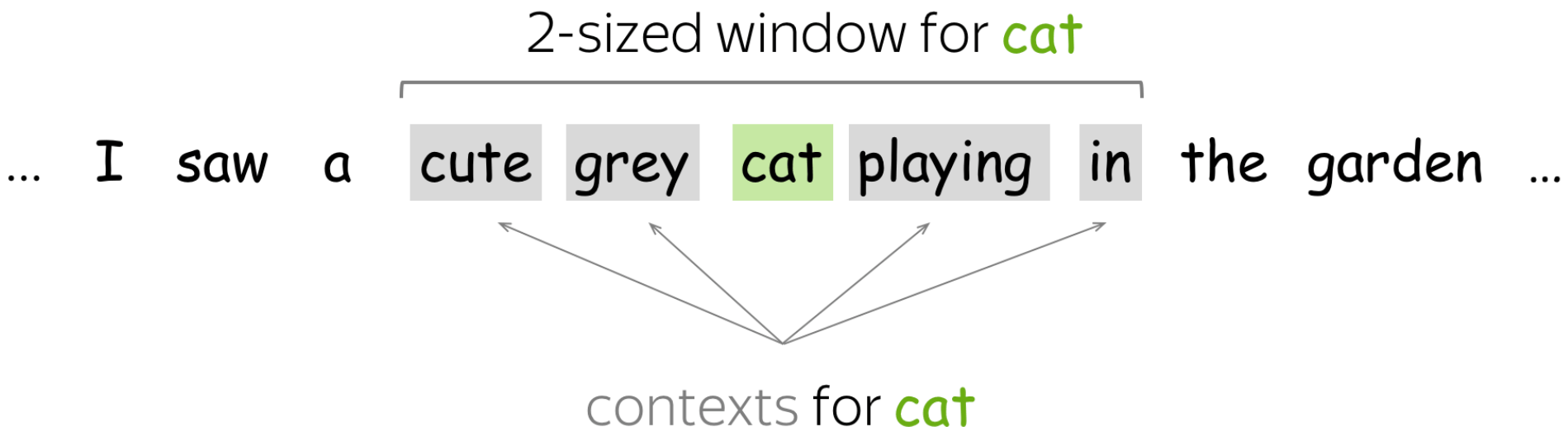


Embedding dimension =
vocabulary size

Issues:

- very high dimensional
- do not capture semantic similarity between words (recall last lecture)

Recap: latent semantic analysis



Context:

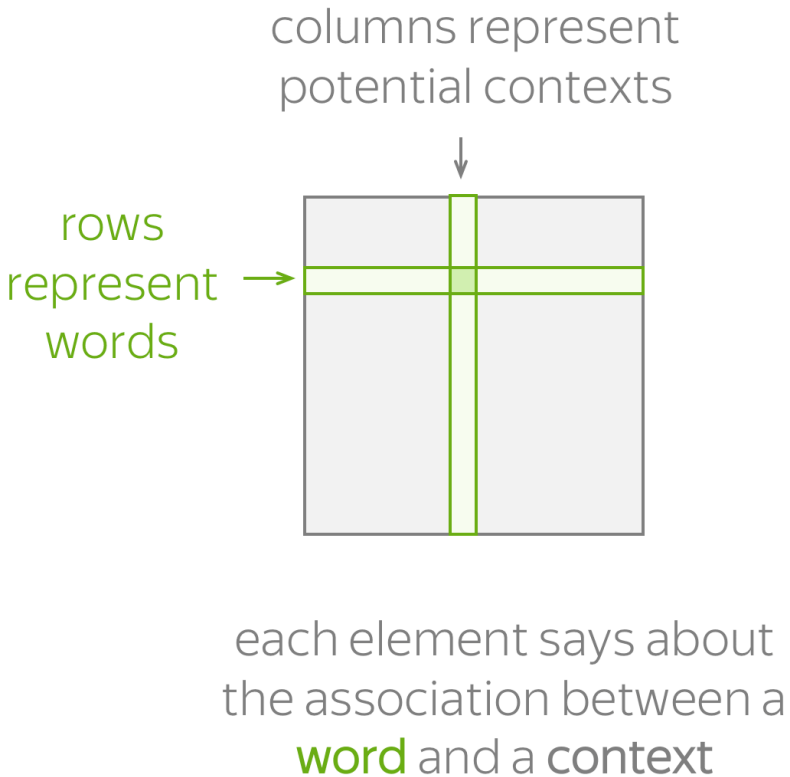
- surrounding words in a L-sized window

Matrix element:

- $N(w, c)$ – number of times **word** w appears in context c

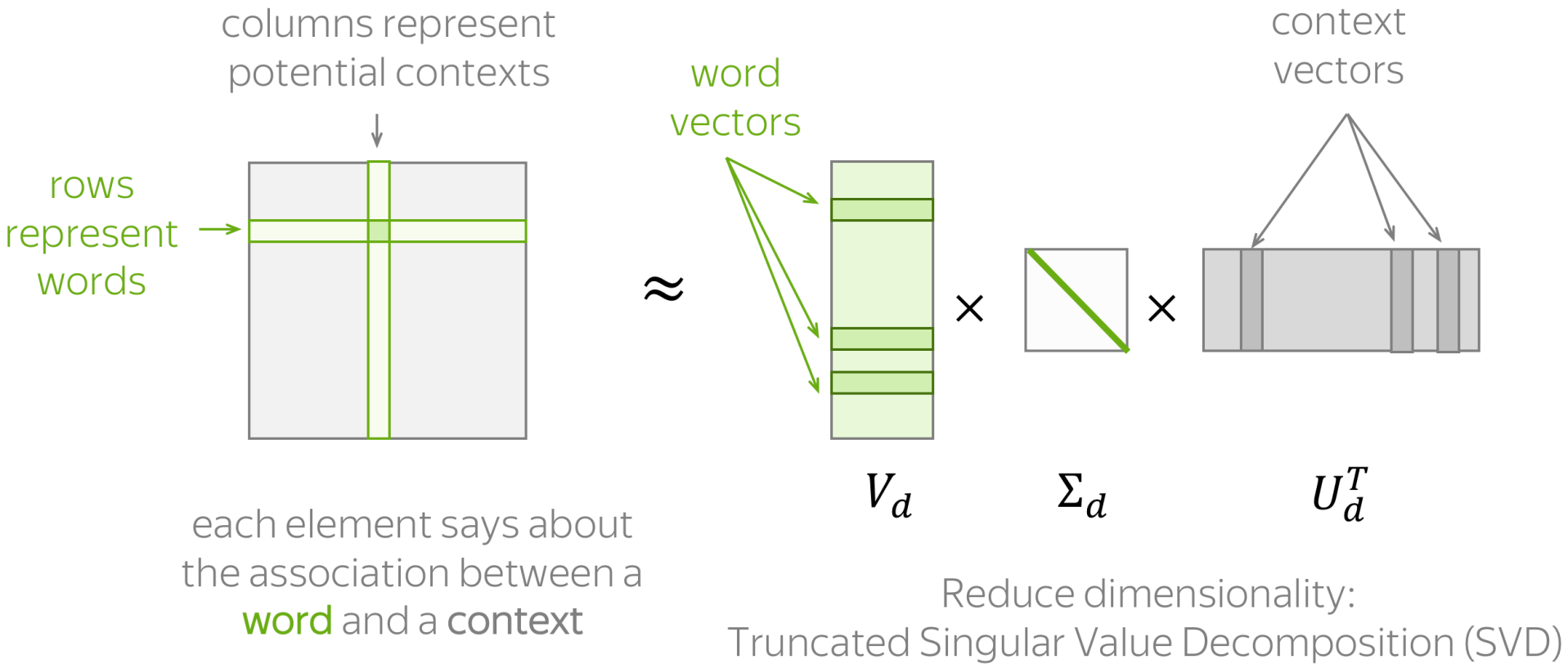
Recall: to make it work reasonably well, you need something more sophisticated (e.g., PMI)

Latent Semantic Analysis



This is either the 'raw' co-occurrence matrix N , or its transformations (e.g., PMI)

Latent Semantic Analysis



Hard to scale, any alternative?

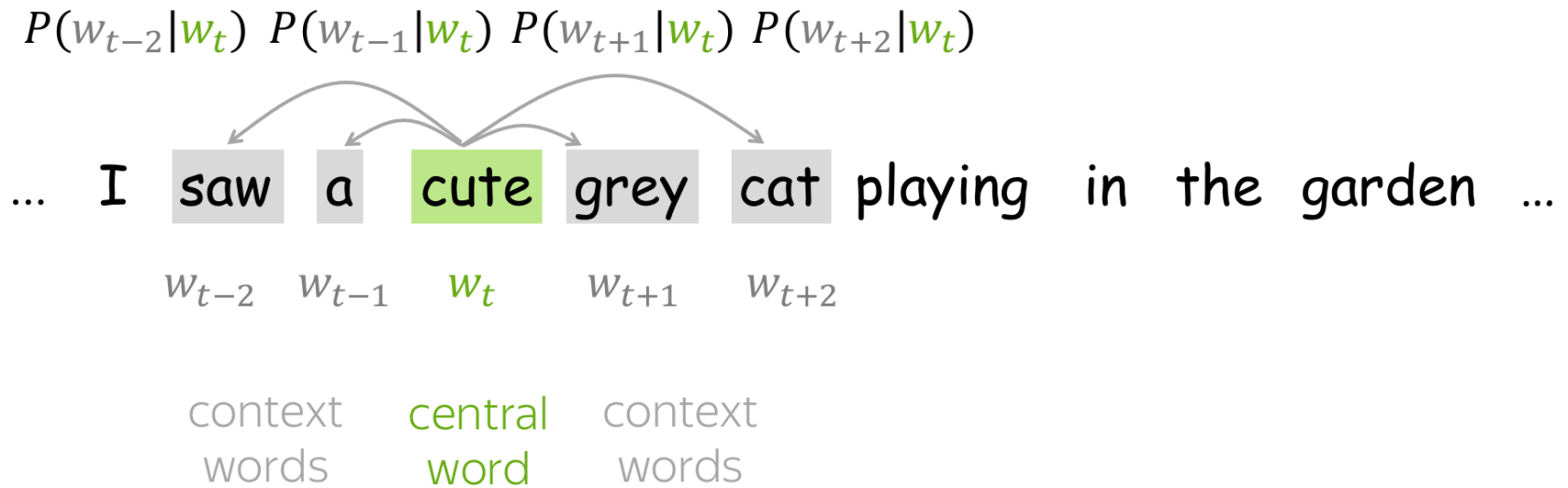
Prediction-based (aka neural) methods

$$P(w_{t-2}|w_t) \quad P(w_{t-1}|w_t) \quad P(w_{t+1}|w_t) \quad P(w_{t+2}|w_t)$$

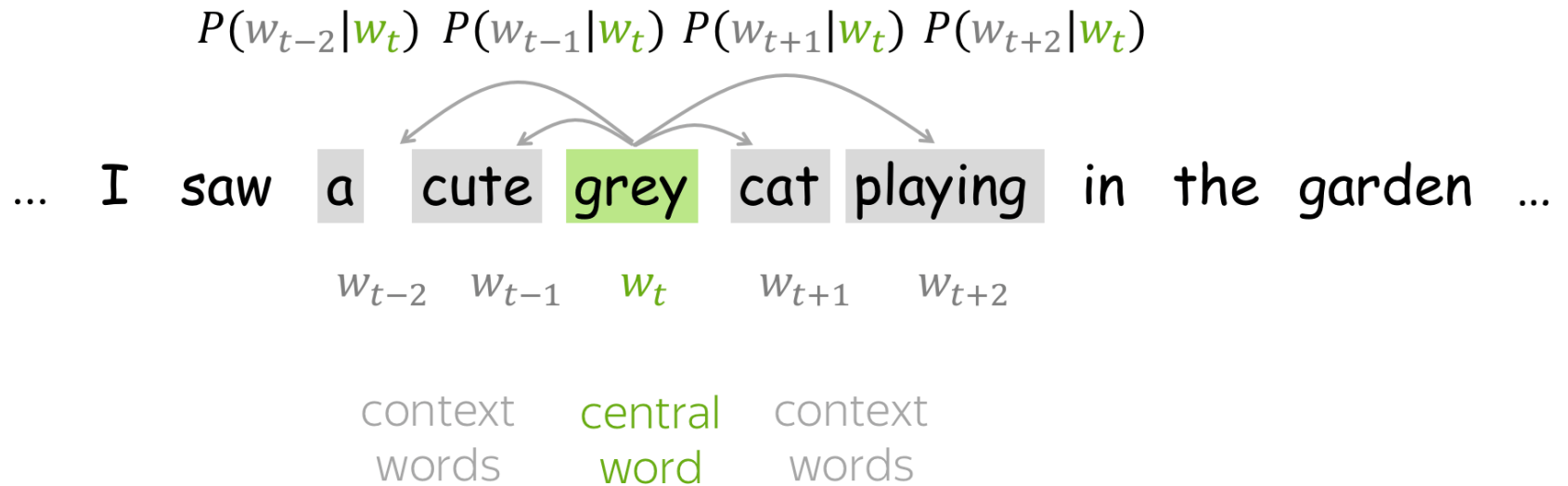


context central context
words word words

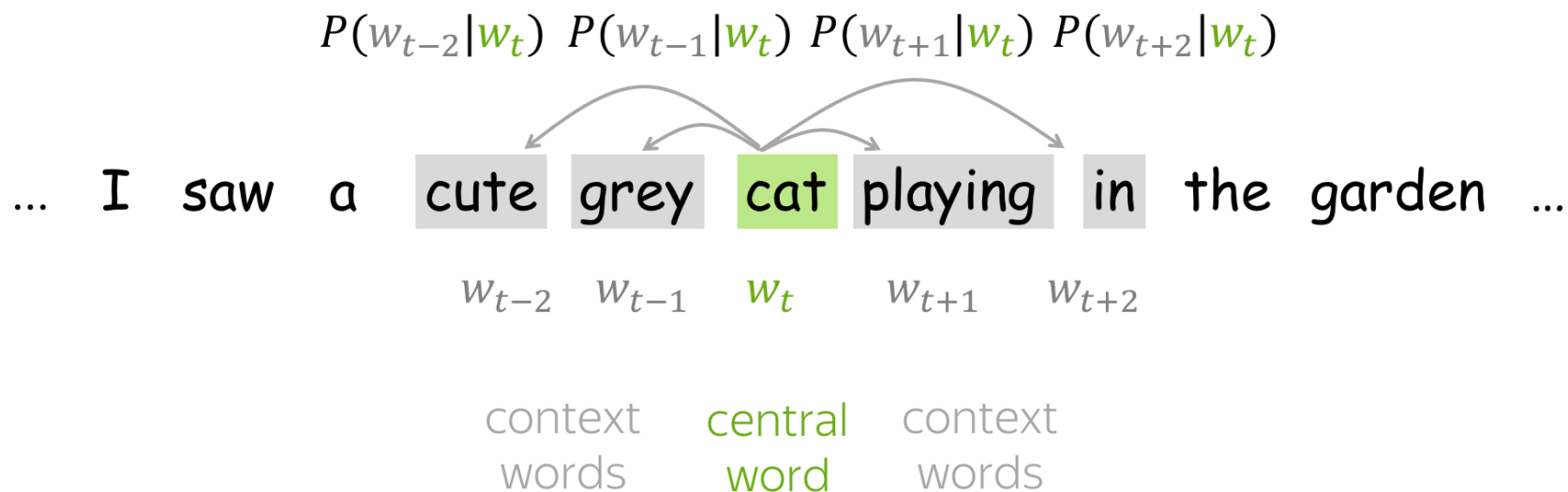
Prediction-based (aka neural) methods



Prediction-based (aka neural) methods

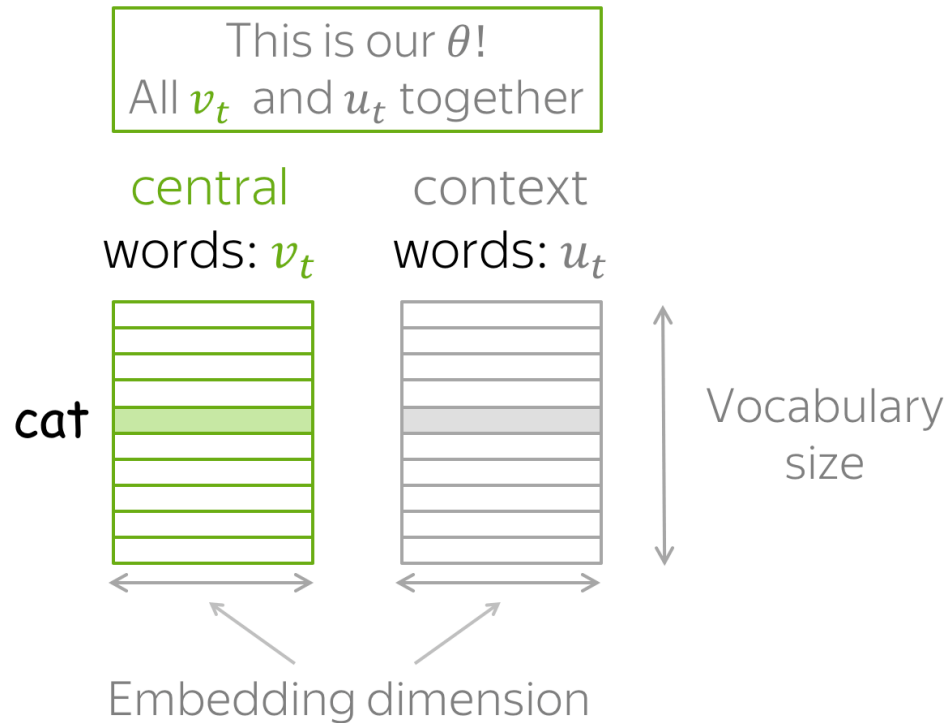


Prediction-based (aka neural) methods



What we discuss is Mikolov's Skipgram but there
several variations on this idea

How do we calculate the probabilities $P(w_{t+j}|w_t, \theta)$?



For each word w we have two vectors:

- when it is a **central word**
- when it is a **context word**

How do we calculate the probabilities $P(w_{t+j}|w_t, \theta)$?

The probability of the context word o given the central word c is

$$P(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)}$$

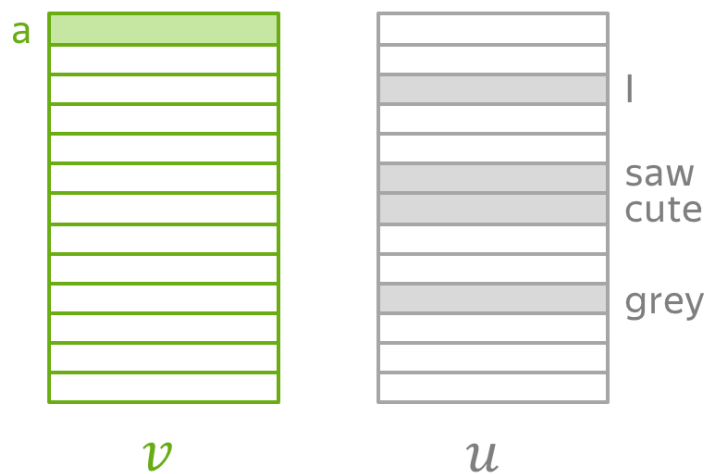
Dot product: measures similarity of o and c
Larger dot product = larger probability

Normalize over entire vocabulary
to get probability distribution

This is the softmax function

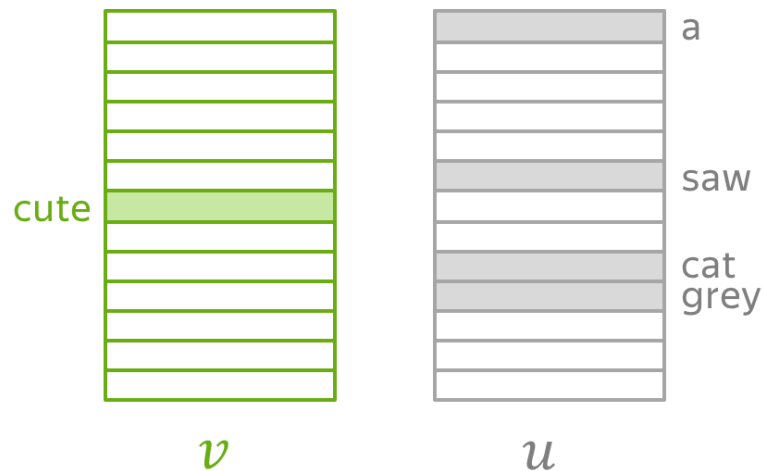
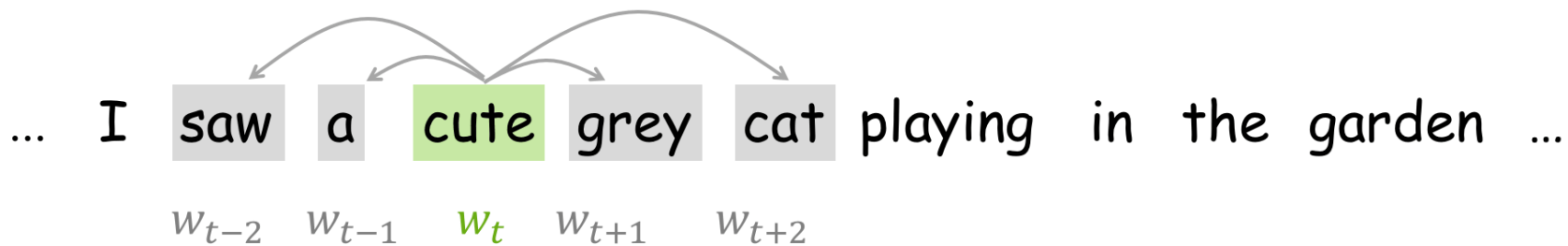
Back to our example

$$P(u_I | v_a) \quad P(u_{saw} | v_a) \quad P(u_{cute} | v_a) \quad P(u_{grey} | v_a)$$



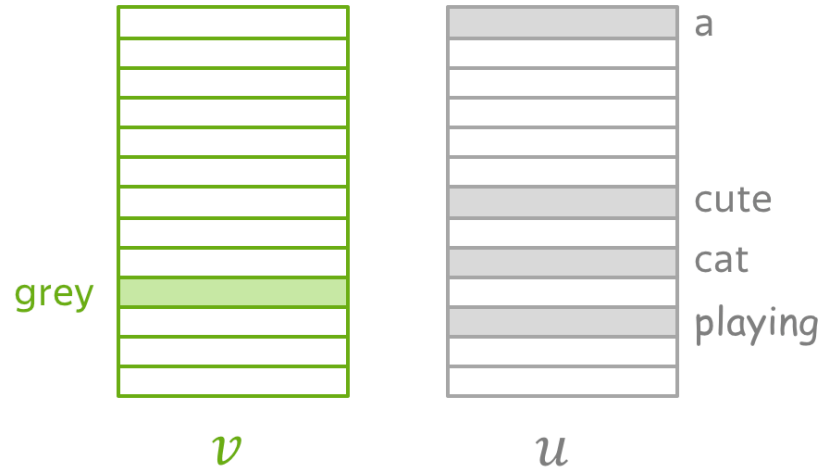
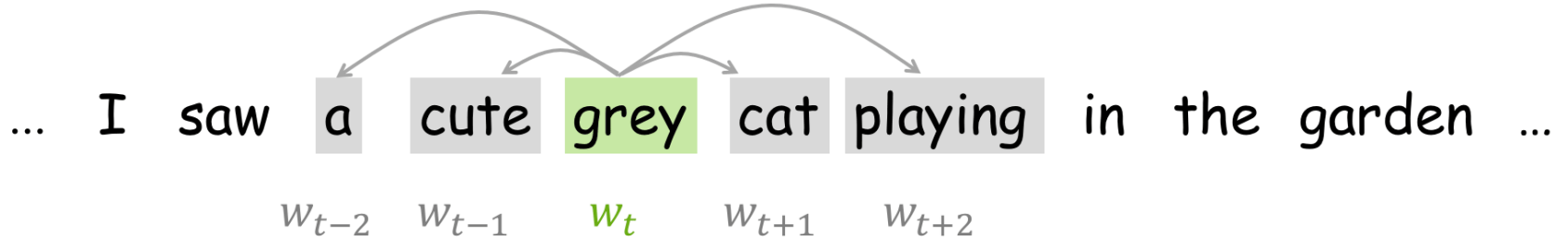
Back to our example

$$P(u_{saw} | v_{cute}) P(u_a | v_{cute}) P(u_{grey} | v_{cute}) P(u_{cat} | v_{cute})$$



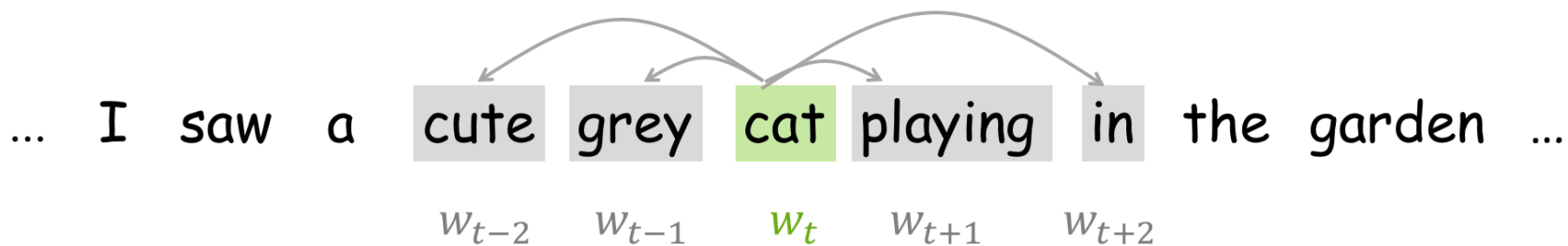
Back to our example

$$P(u_a | v_{grey}) P(u_{cute} | v_{grey}) P(u_{cat} | v_{grey}) P(u_{playing} | v_{grey})$$



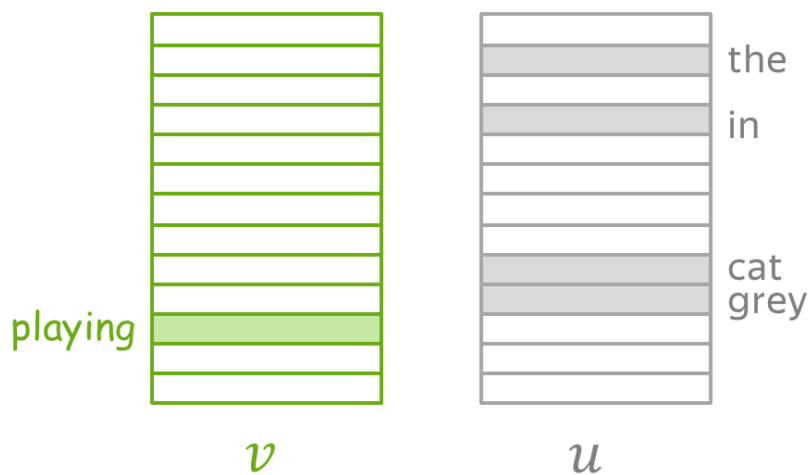
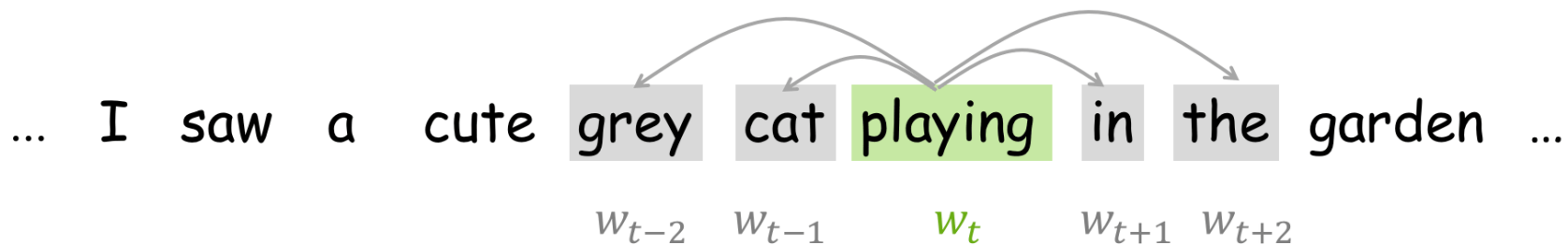
Back to our example

$$P(u_{cute}|v_{cat}) \quad P(u_{grey}|v_{cat}) \quad P(u_{playing}|v_{cat}) \quad P(u_{in}|v_{cat})$$



Back to our example

$$P(u_{\text{grey}}|v_{\text{playing}}) \quad P(u_{\text{cat}}|v_{\text{playing}}) \quad P(u_{\text{in}}|v_{\text{playing}}) \quad P(u_{\text{the}}|v_{\text{playing}})$$



Back to our example

$$P(u_{grey}|v_{in}) \quad P(u_{cat}|v_{in}) \quad P(u_{in}|v_{in}) \quad P(u_{the}|v_{in})$$

... I saw a cute grey **cat** **playing** **in** **the** **garden** ...

w_{t-2} w_{t-1} w_t w_{t+1} w_{t+2}



v



u

What do we optimize?

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

agrees with our
plan above



go over text



with a sliding
window



compute probability of the
context word given the central



How do we optimize the model?

We rely on gradient descent (recall the Logistic Regression lecture)

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta).$$

In practice, we optimize **one word at a time**:

$$-\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} | w_t, \theta) = \frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} J_{t,j}(\theta)$$

How do we optimize the model?

... I saw a cute grey cat playing in the garden ...

$$J_{t,j}(\theta) = -\log P(\text{cute}|\text{cat}) = -\log \frac{\exp u_{\text{cute}}^T v_{\text{cat}}}{\sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}}$$

How do we optimize the model?

... I saw a cute grey cat playing in the garden ...

$$J_{t,j}(\theta) = -\log P(\text{cute}|\text{cat}) = -\log \frac{\exp u_{\text{cute}}^T v_{\text{cat}}}{\sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}}$$

$$= -u_{\text{cute}}^T v_{\text{cat}} + \log \sum_{w \in \text{Voc}} \exp u_w^T v_{\text{cat}}$$

Note which parameters are present at this step:

- from vectors for central words, only v_{cat} ;
- from vectors for context words, all u_w (for all words in the vocabulary)

How do we optimize the model?

... I saw a cute grey cat playing in the garden ...

$$J_{t,j}(\theta) = \underbrace{-u_{cute}^T v_{cat}}_{\textcircled{1}} + \log \underbrace{\sum_{w \in V} \exp(u_w^T v_{cat})}_{\textcircled{2}}$$

make an update

$$v_{cat} := v_{cat} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{cat}}$$

$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$

How do we optimize the model?

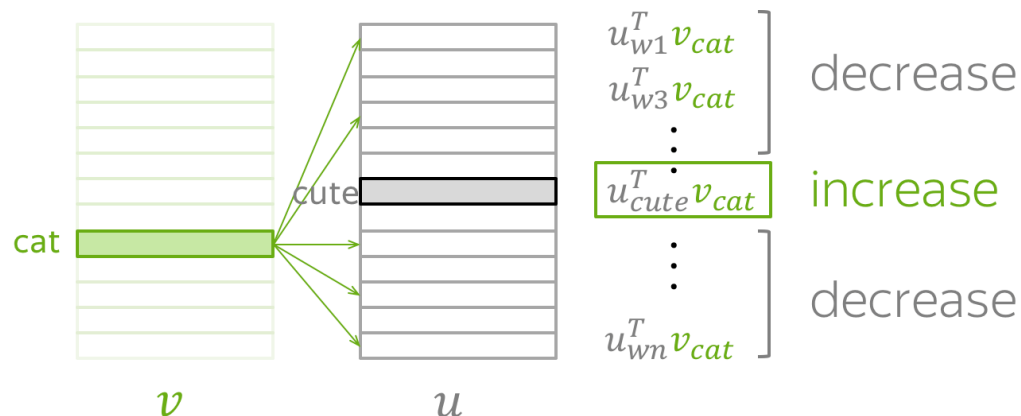
... I saw a cute grey cat playing in the garden ...

$$J_{t,j}(\theta) = \underbrace{-u_{cute}^T v_{cat}}_{\textcircled{1}} + \log \underbrace{\sum_{w \in V} \exp(u_w^T v_{cat})}_{\textcircled{2}}$$

make an update

$$v_{cat} := v_{cat} - \alpha \frac{\partial J_{t,j}(\theta)}{\partial v_{cat}}$$

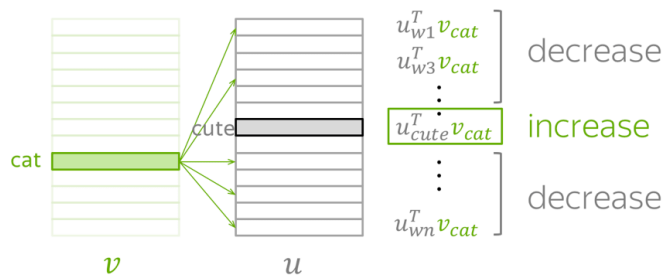
$$u_w := u_w - \alpha \frac{\partial J_{t,j}(\theta)}{\partial u_w} \quad \forall w \in V$$



Making learning more efficient

Dot product of v_{cat} :

- with u_{cute} - increase,
- with all other u - decrease



Parameters to be updated:

- v_{cat}
- u_w for all w in the vocabulary $|V| + 1$ vectors

Making learning more efficient

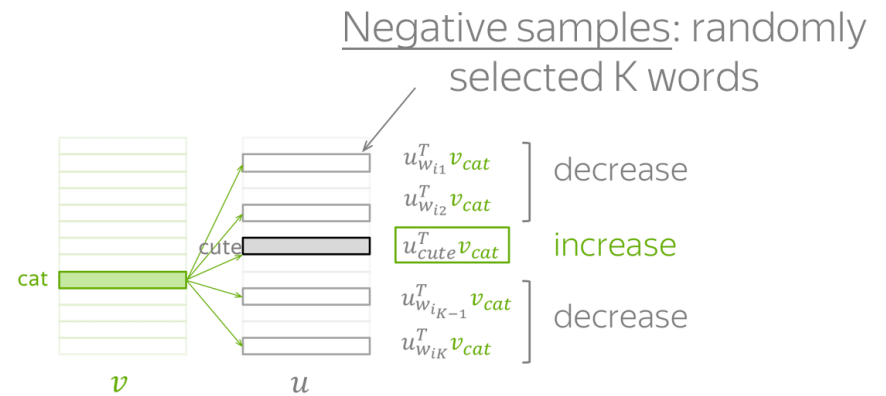
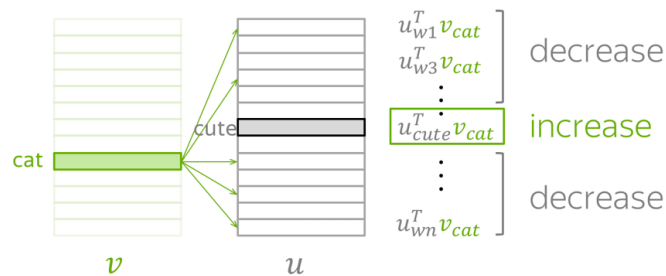
Dot product of v_{cat} :

- with u_{cute} - increase,
- with all other u - decrease



Dot product of v_{cat} :

- with u_{cute} - increase,
- with a subset of other u - decrease



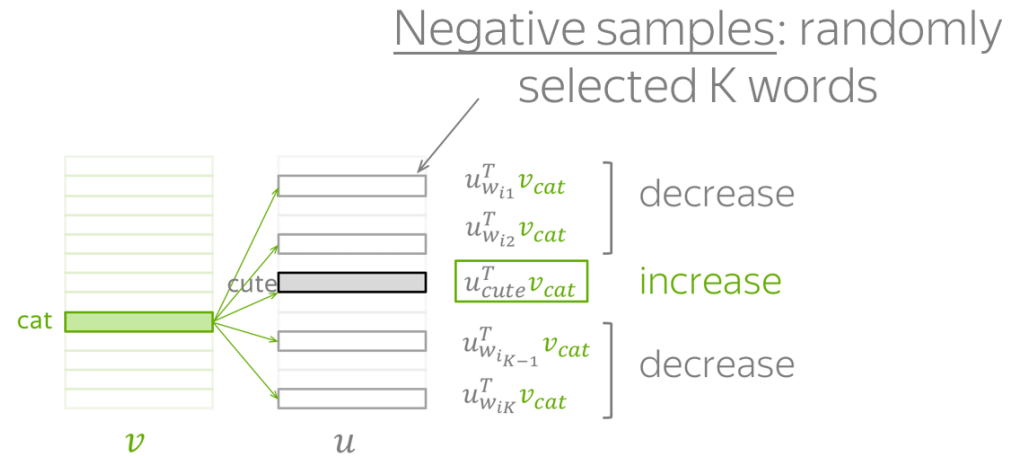
Parameters to be updated:

- v_{cat}
- u_w for all w in the vocabulary $|V| + 1$ vectors

Parameters to be updated:

- v_{cat}
- u_{cute} and u_w for w in K negative examples $K + 2$ vectors

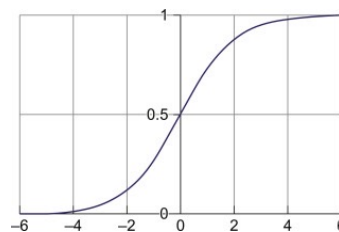
Negative sampling



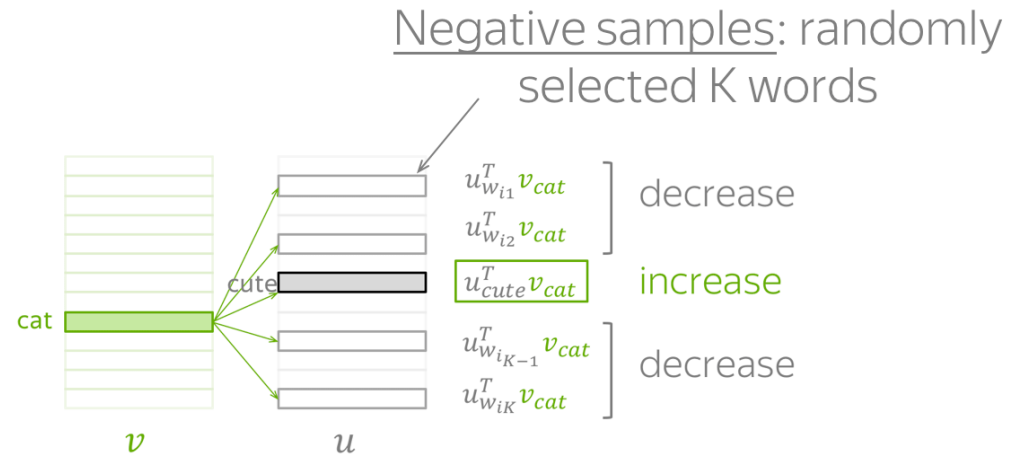
$$J_{t,j}(\theta) = -\log \sigma(u_{cute}^T v_{cat}) - \sum_{w \in \{w_{i1}, \dots, w_{iK}\}} \log(1 - \sigma(u_w^T v_{cat}))$$

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

the logistic sigmoid function



Negative sampling



Converted to binary classification: predict "+" for (central, real context) pairs, and "-" for (central, random context)

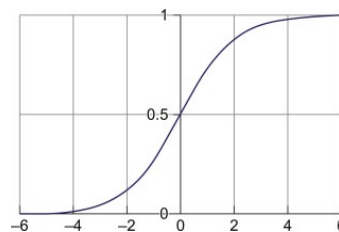
“Positive examples”

“Negative examples”

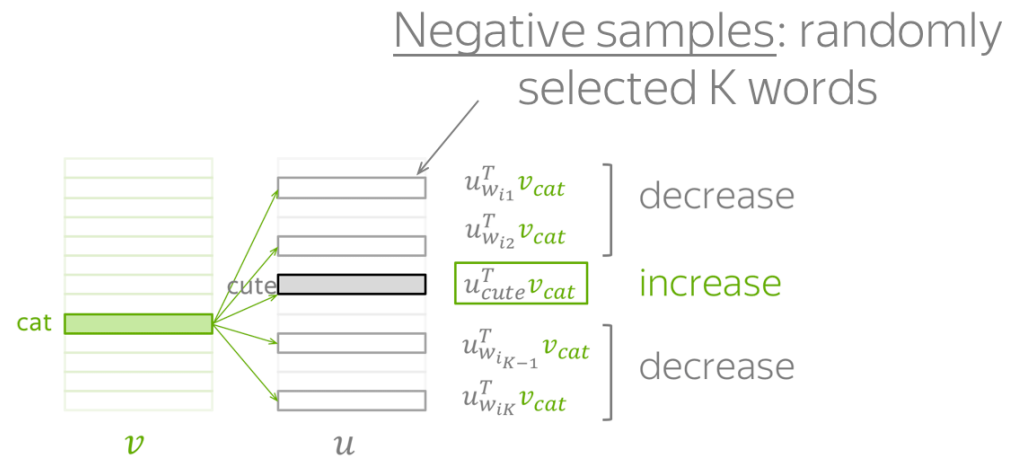
$$J_{t,j}(\theta) = -\log \sigma(u_{cute}^T v_{cat}) - \sum_{w \in \{w_{i1}, \dots, w_{iK}\}} \log(1 - \sigma(u_w^T v_{cat}))$$

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

the logistic sigmoid function

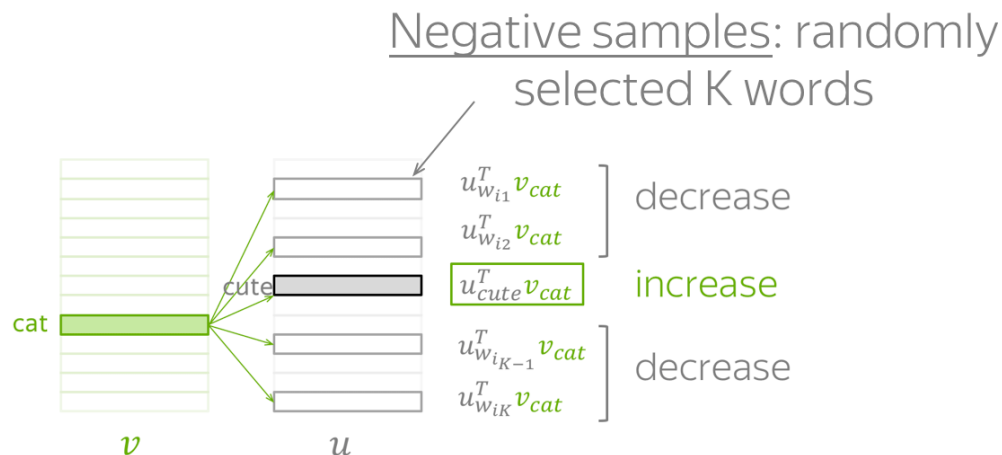


How to select negative samples?



The basic idea is to select random words based on their (unigram) frequency in a corpus, what are the issues and how can this be improved?

How to select negative samples?



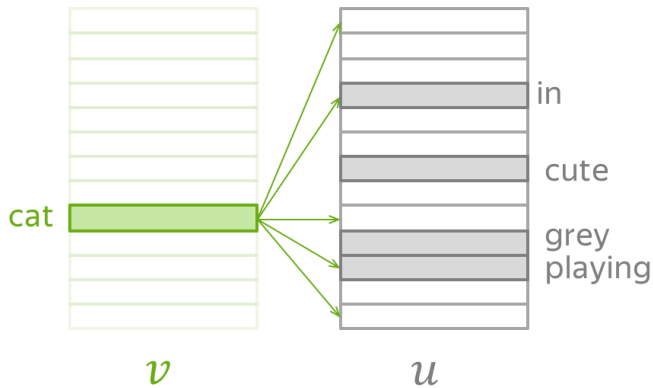
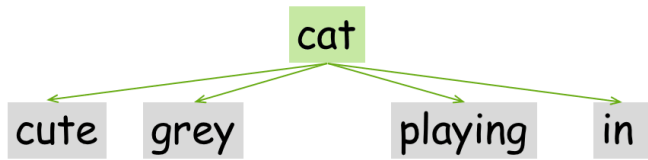
The basic idea is to select random words based on their (unigram) frequency in a corpus, what are the issues and how can this be improved?

- “Flatten” the unigram distribution to make sure infrequent words get sampled (recall Zipf’s distribution)
- Don’t generate compatible contexts
- Make sure you generate “hard” negative examples, to learn informative word representations

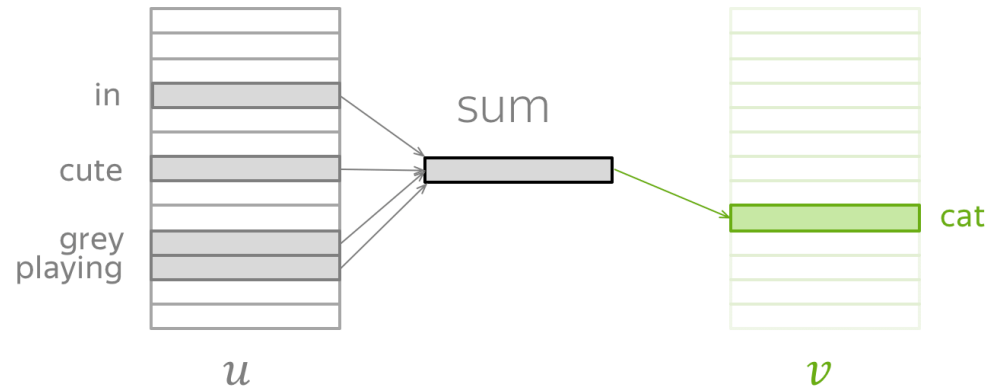
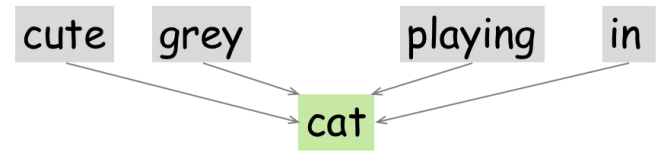
Negative sampling is an important technique used in many contexts in NLP/ML, often with these mods

Variations of Word2Vec

... I saw a cute grey cat playing in the garden ...



Skip-Gram: from central predict context
(one at a time)



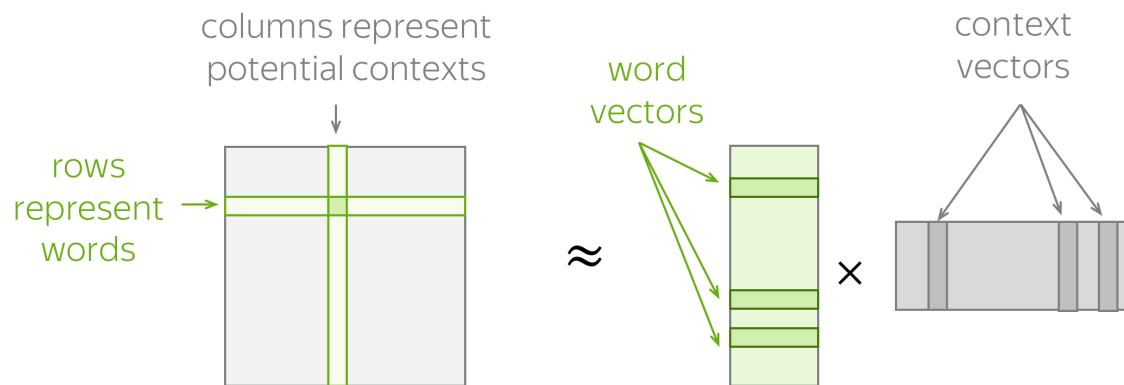
CBOW: from sum of context predict central

Generally, word2vec was far from the first idea of learning embedding, its success is largely due to simplicity (+ the efficient implementation and stability)

Relation to LSA

It is possible to show that optimizing the skipgram objective (with the negative sampling modification) corresponds to factorizing PMI matrix (Levy & Goldberg 2014)

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t, \theta)$$



Evaluation (~ recap)

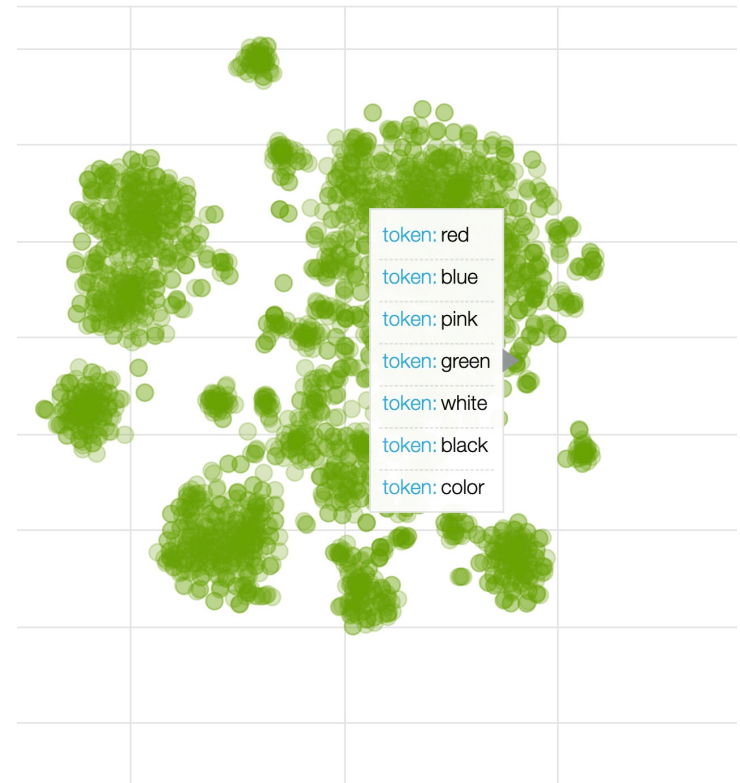
Intrinsic

Use embeddings to represent tokens within models for downstream tasks (e.g., sentiment classification, question answering, ...)

Extrinsic

- Relatedness
- Word associations
- **Analogy**

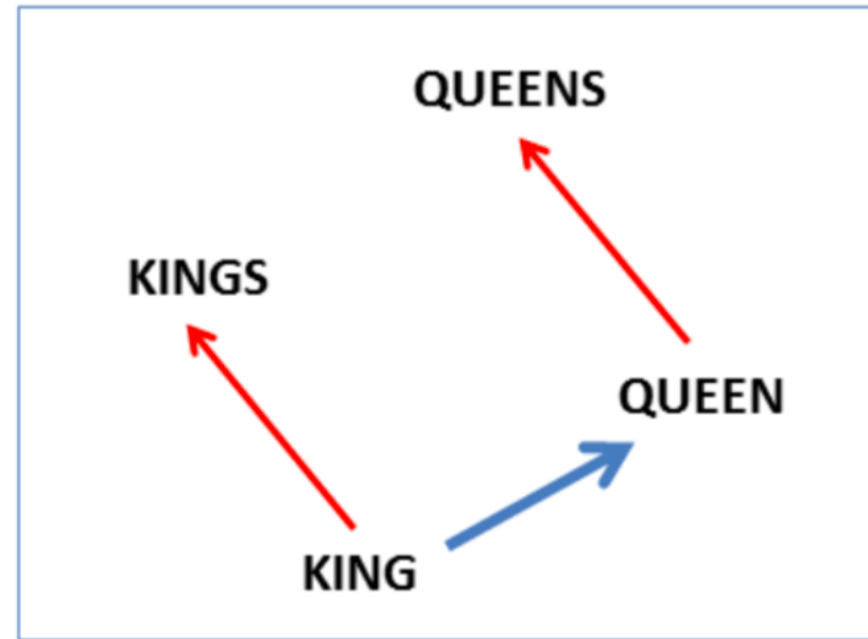
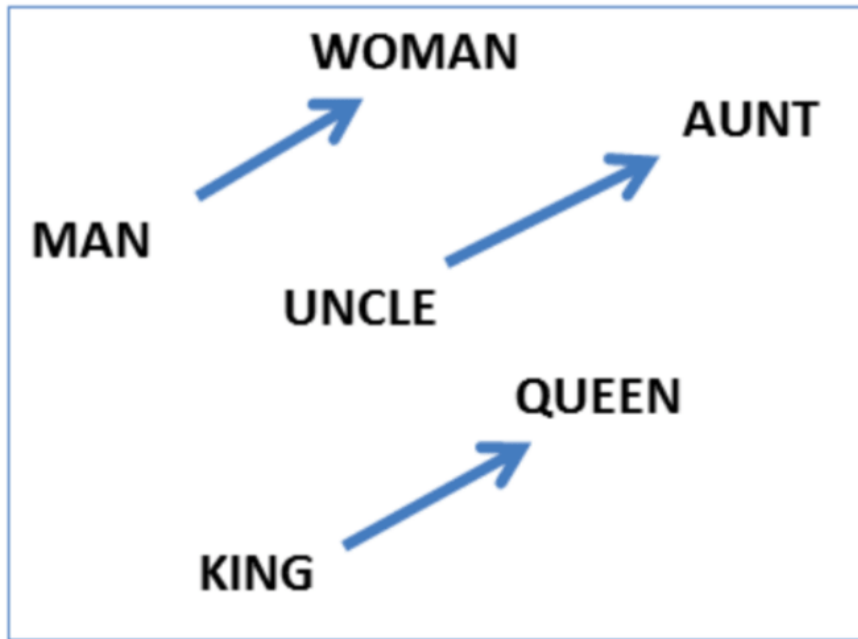
As we will see soon the low-dimensional vectors produced by neural methods are much 'friendlier' to downstream applications than sparse counts



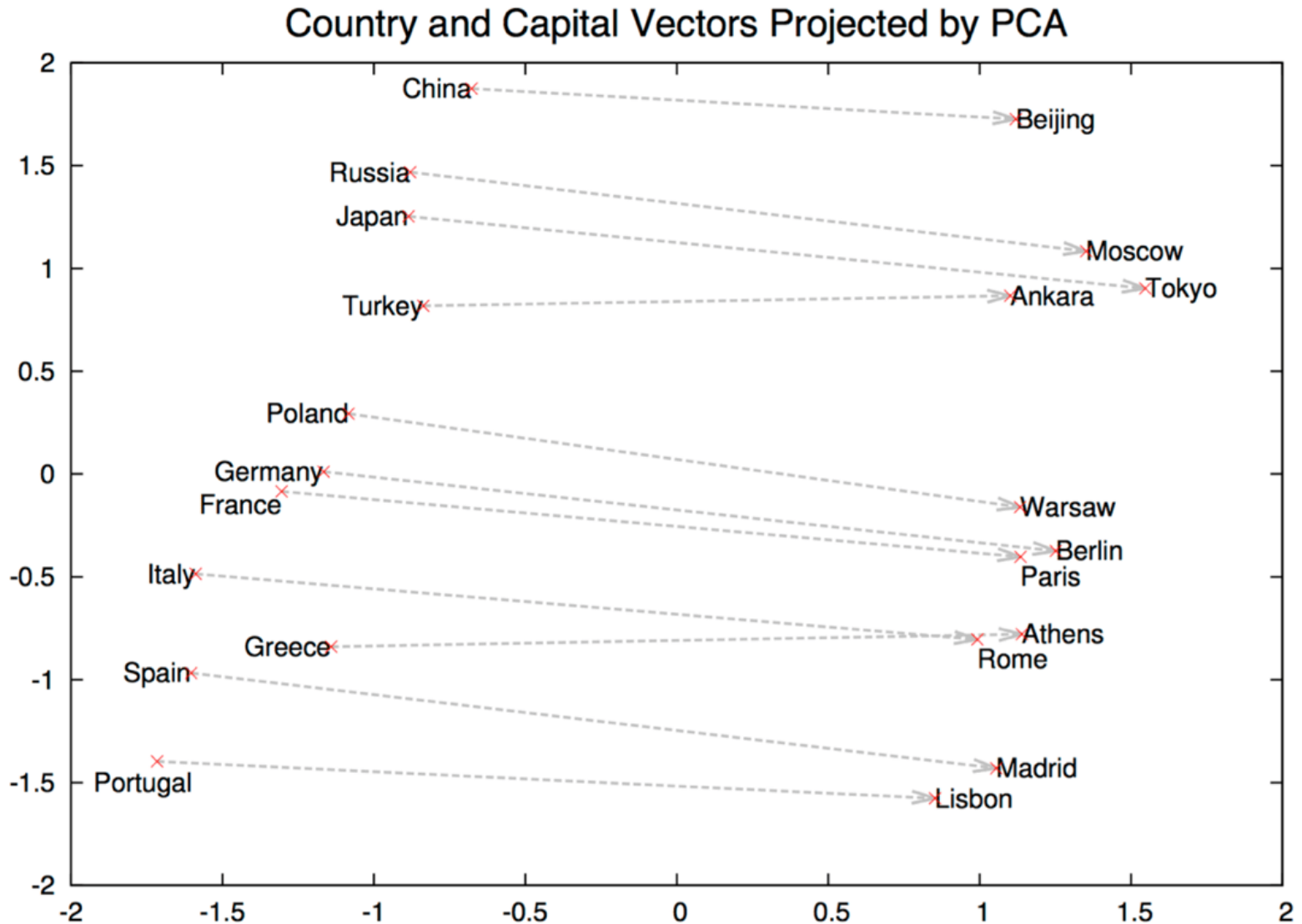
Analogy

semantic: $v(\text{king}) - v(\text{man}) + v(\text{woman}) \approx v(\text{queen})$

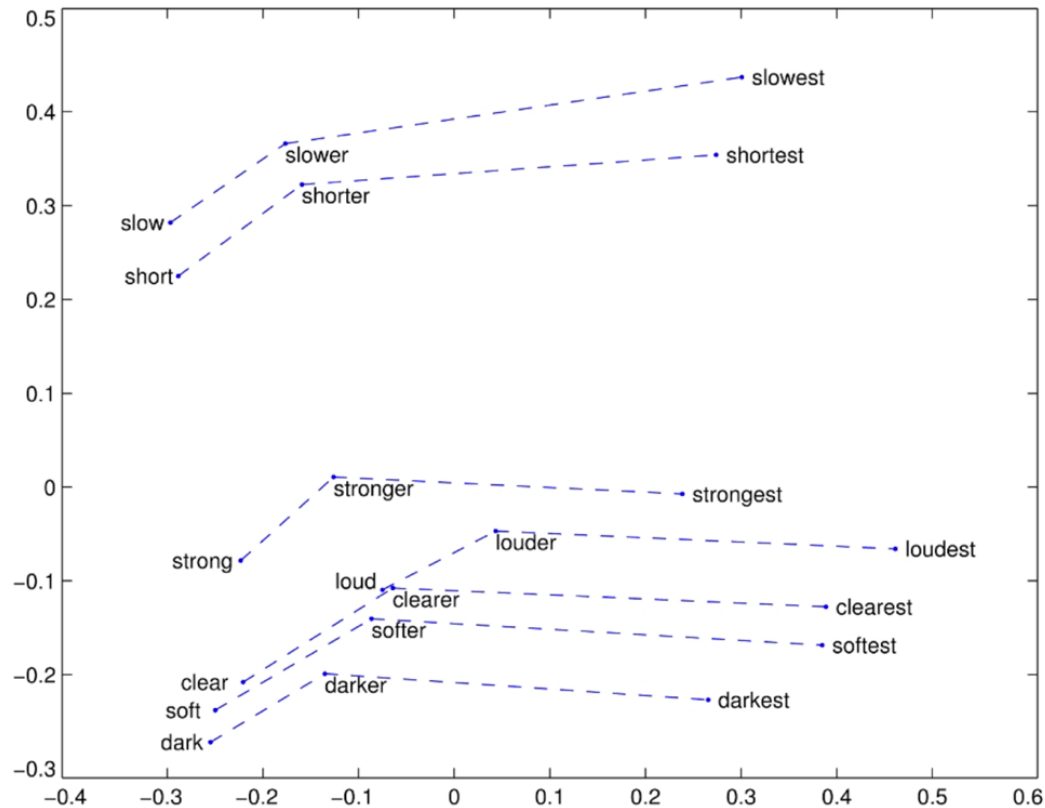
syntactic: $v(\text{kings}) - v(\text{king}) + v(\text{queen}) \approx v(\text{queens})$



Analogy: semantic relations



Analogy: syntactic relations



Why makes word2vec represent relations in such linear way? See a paper by an Edinburgh student: [Analogies Explained: Towards Understanding Word Embeddings](#), Allen & Hospidales, ICML 2019

There are some important caveats about the analogy benchmarks and visualizations: <https://aclanthology.org/N18-2039/>

Summary for today

Neural embeddings can be **efficiently** learned from large collections of unannotated texts ('self-supervision')

Many algorithms, and important **hyperparameter** choices (windows sizes, numbers of negatives samples)

Useful in practice and have some intriguing properties

Preferable over raw count-based method but:

- how do we handle **multiple senses?** (ambiguity)
- how do we encode **longer spans of text?** (compositionality)

Check out Lena Voita's online resource – NLP class for you:

https://lena-voita.github.io/nlp_course.html