# Autoencoders
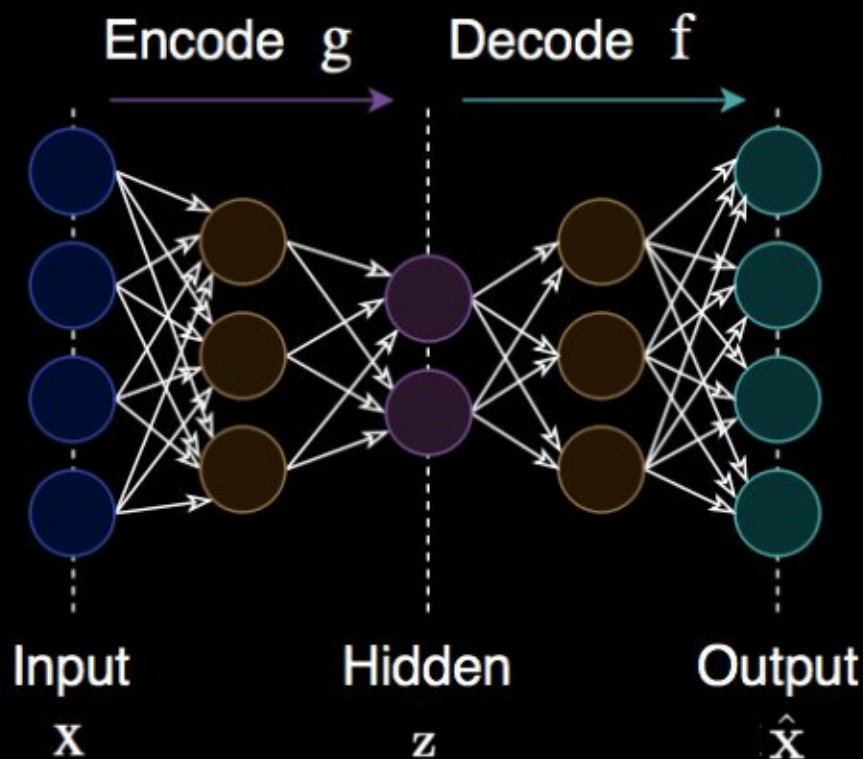
Denote **z** as encoded with encoder E input **x**

$$\mathbf{z} = E(\mathbf{x}, \boldsymbol{\theta}_E)$$

Decoder D recovers **x** from latent representation

$$\hat{\mathbf{x}} = D(\mathbf{z}, \boldsymbol{\theta}_D)$$

Optimal parameters learned w.r.t. loss function L

$$[\boldsymbol{\theta}_E, \boldsymbol{\theta}_D] = \arg\min L(\hat{\mathbf{x}}, \mathbf{x})$$

Encode $g$  Decode $f$



Input
**x**

Hidden
**z**

Output
$\hat{\mathbf{x}}$

Image source:

Denote **z** as encoded with encoder E input **x**

$$\mathbf{z} = E(\mathbf{x}, \boldsymbol{\theta}_E)$$

Decoder D recovers **x** from latent representation

$$\hat{\mathbf{x}} = D(\mathbf{z}, \boldsymbol{\theta}_D)$$

Simple example: PCA

Optimal parameters learned w.r.t. loss function L

$$[\boldsymbol{\theta}_E, \boldsymbol{\theta}_D] = \arg\min L(\hat{\mathbf{x}}, \mathbf{x})$$

Image source: Habr post on autoencoders and GANs

16 components

# Convolutional performance on MNIST



7 x 7 latent space

- **One-hot vectors:**

```
Rome  = [1, 0, 0, 0, 0, 0, …, 0]

Paris = [0, 1, 0, 0, 0, 0, …, 0]

Italy = [0, 0, 1, 0, 0, 0, …, 0]

France = [0, 0, 0, 1, 0, 0, …, 0]
```

**Problems:**

- Huge vectors
- VERY sparse
- No semantics or word similarity information included

# Why not to learn word vectors?

What is `king` - `man` + `woman`?

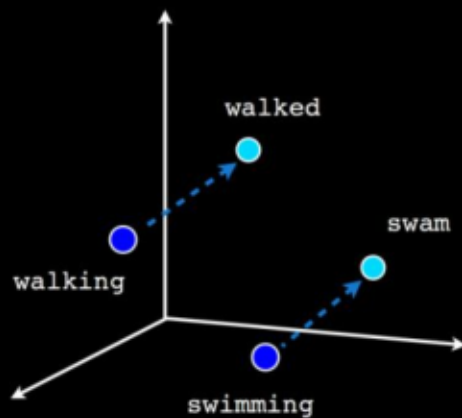So `king` - `man` + `woman` = `queen`!

- **Word2vec** (Mikolov et al. 2013) - a framework for learning word embeddings



Male-Female

Verb tense

Country-Capital
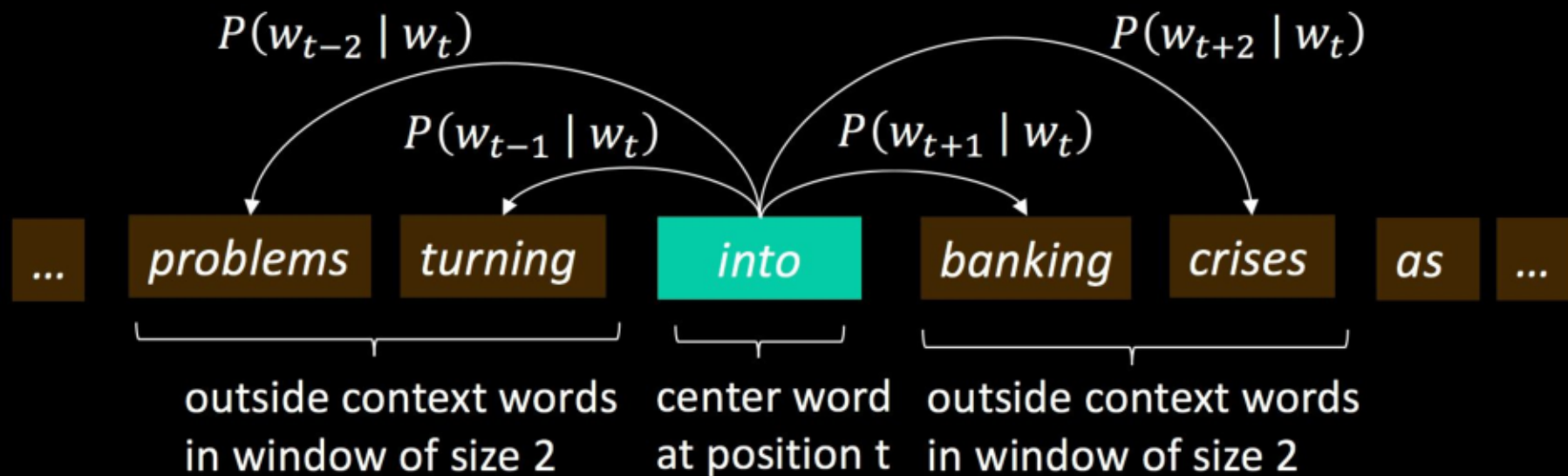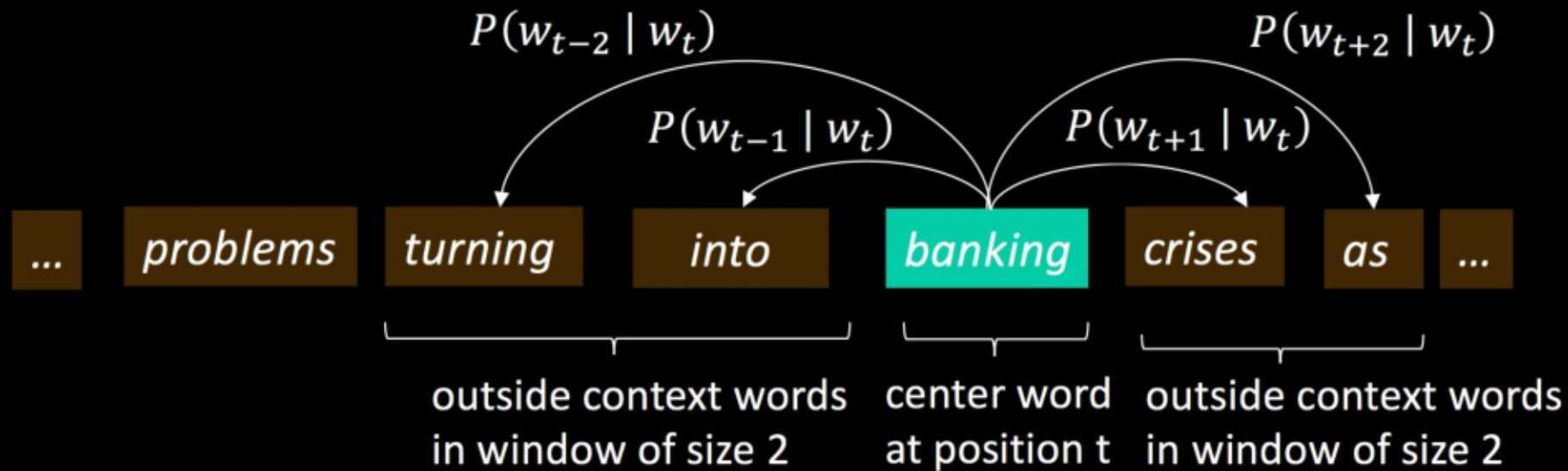
# Embeddings: word2vec



Source Text

Training Samples

The **quick** brown fox jumps over the lazy dog. ⟹ (the, quick) (the, brown)

The **quick** brown fox jumps over the lazy dog. ⟹ (quick, the) (quick, brown) (quick, fox)

The quick **brown** fox jumps over the lazy dog. ⟹ (brown, the) (brown, quick) (brown, fox) (brown, jumps)

The quick brown **fox** jumps over the lazy dog. ⟹ (fox, quick) (fox, brown) (fox, jumps) (fox, over)

45

$P(w_{t-2} \mid w_t)$

$P(w_{t+2} \mid w_t)$

$P(w_{t-1} \mid w_t)$

$P(w_{t+1} \mid w_t)$

... problems turning into banking crises as ...

outside context words
in window of size 2

center word
at position t

outside context words
in window of size 2

Source: CS224n: http://web.stanford.edu/class/cs224n/slides/cs224n-2019-lecture01-wordvecs1.pdf

# Embeddings: word2vec



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

... problems turning into banking crises as ...

outside context words in window of size 2

center word at position t

outside context words in window of size 2

# Embeddings: word2vec

Source: http://mccormickml.com/2016/04/19/word2vec-tutorial-the-skip-gram-model/

# GloVe Visualizations
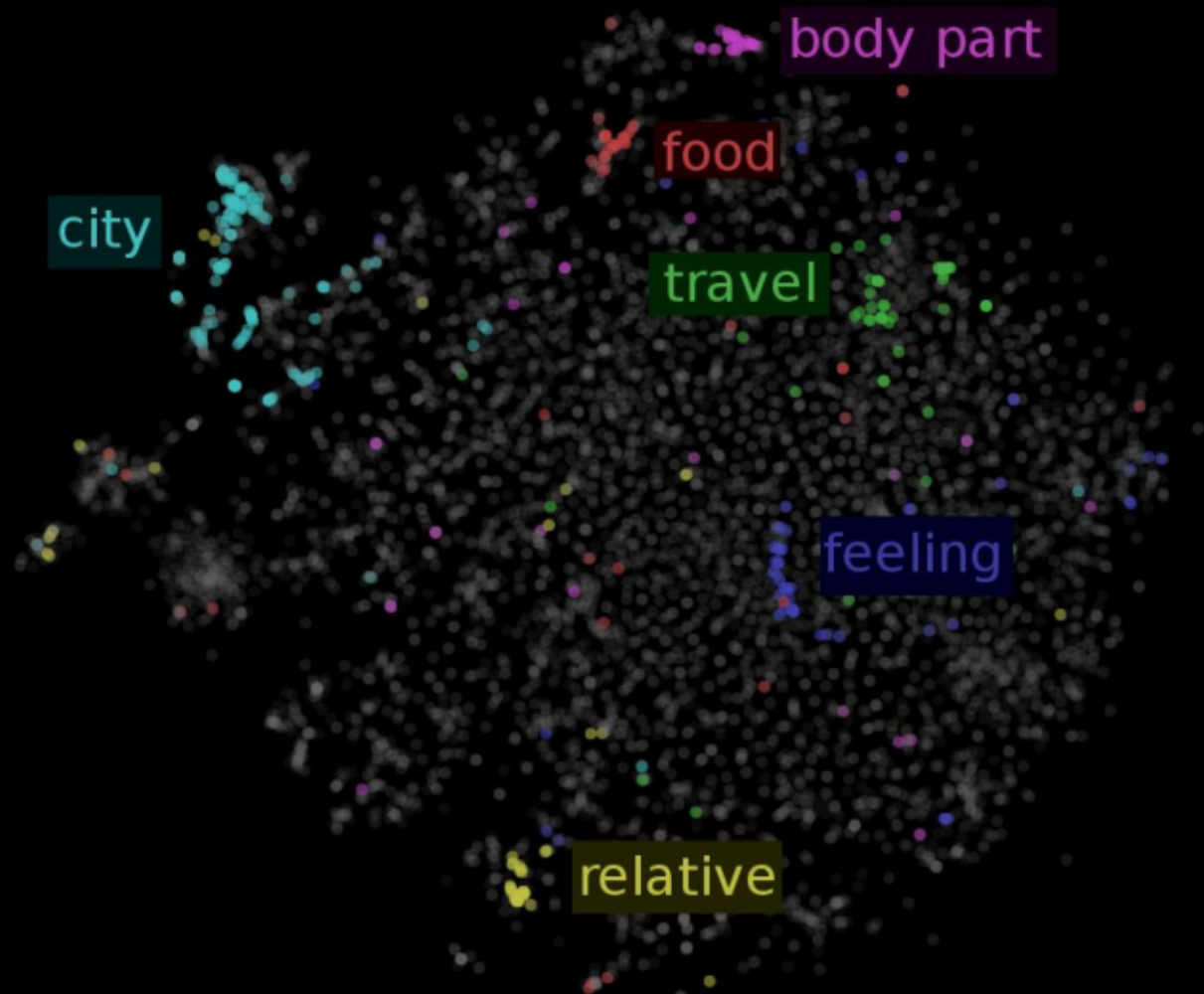
King - man + woman = queen

$x$      $y$      $y'$      $target$

$$\cos(x-y+y', target) \rightarrow \max_{target}$$

60

- **Use statistics:**
  - T-criterion

$$t = \frac{\overline{x} - \mu}{\sqrt{\dfrac{s^2}{N}}}$$

$H_0$ : 'social media' occurs with probability:

$$\mu = P(social)P(media) = \frac{C(social)(media)}{N^2}$$

$H_a$ : 'social media' does not occur with such a probability

- **Use statistics:**
  - Chi-squared

$$\chi^2 = \sum_{ij} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$E(social\ media) = \frac{C(social)}{N} \cdot \frac{C(media)}{N} \cdot N$$

$$O_{ij}\ from\ table$$

|  | w1 = social | w1 != social |
|---|---|---|
| **w2 = media** | C(social media) | C(x media) where x could be any word |
| **w2 != media** | C(social x) where x could be any word | C(any pair not starting with social or ending with media) |