

# Anomaly Detection in Probabilistic Databases with Embeddings of Categorical Values

---

Daoping Wang

March 20, 2020

Technische Universität München  
Mercateo Deutschland AG

1. Introduction
2. Probabilistic Database Modeling
3. Learning Embeddings
4. Experimental Results

- Enabler for advanced analytics
- Key resource for ML-based techniques
- Requirement on data quantity and **quality**
  - Accuracy, consistency and completeness
  - Ground truth information



---

<sup>1</sup><https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>

- Enabler for advanced analytics
- Key resource for ML-based techniques
- Requirement on data quantity and **quality**
  - Accuracy, consistency and completeness
  - Ground truth information



How does real-world data look like?

---

<sup>1</sup><https://www.economist.com/leaders/2017/05/06/the-worlds-most-valuable-resource-is-no-longer-oil-but-data>

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	date
102-805270-BP	1039U	9739011201	Knipex	Eckschwindenzange	C62	55.7	2022-06-22
102-805270-BP	1082	1520050	Knipex	Inconsistency	C62	14.82	4003773022022
102-805270-BP	108EL	Conflicts	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Premium	Ratschenkabelschere	C62	141.69	4003773043935

Outlier (?)

Missing values

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	date
102-805270-BP	1039U	9739011201	Knipex	Eckschwendenzange	C62	55.7	2022-06-22
102-805270-BP	1082	1520050	Knipex	Inconsistency	C62	14.82	4003773022022
102-805270-BP	108EL	Conflicts	Knipex	Eck-Rohrzange	C62	18.33	Missing values
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	Missing values
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

Outlier (?)

Missing values

**Goal: Automated anomaly detection!**

- Heuristic methods — *if ... then ...*
  - Maintainability
  - Conflicting rules
  - “Knipex” vs “Primum”: Who is wrong?
  - High demands on domain knowledge
- Statistical methods
  - **High-cardinality:** Most values are rarely observed
  - Occurrence information alone is not enough

- Heuristic methods — *if ... then ...*
  - Maintainability
  - Conflicting rules
  - “Knipex” vs “Primum”: Who is wrong?
  - High demands on domain knowledge
- Statistical methods
  - **High-cardinality:** Most values are rarely observed
  - Occurrence information alone is not enough

**How does a human perform anomaly detection?**

Outlier (?)

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

Outlier (?)

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

### Domain knowledge:

- $P(ek\_amount | keyword, manufacturer) \sim$  “Similar products have similar prices.”
- $P(ek\_amount | set\_id) \sim$  “Articles from the same set have similar prices.”
- $P(ek\_amount | manufacturer) \sim$  “Knipex is cheaper than Wera.”
- ...

A diagram illustrating probabilistic relationships between table columns. Red ovals highlight the columns: set\_id, catalog\_id, manufacturer, keyword, unit, and ek\_amount. Red curved arrows show dependencies: catalog\_id points to manufacturer, manufacturer points to keyword, and unit points to ek\_amount.

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschweidenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

### Accumulating domain knowledge:

- Modeling joint probability distribution

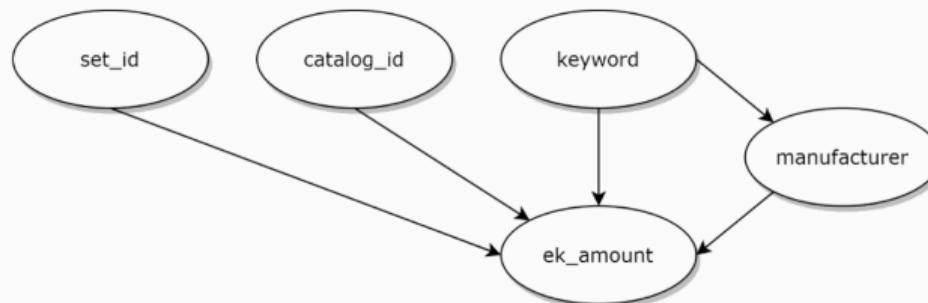
$$P(\text{set\_id}, \text{catalog\_id}, \text{manufacturer}, \text{keyword}, \text{ek\_amount})$$

- Marginal probability = Confidence score

<i>set_id</i>	$P(\text{set\_id})$
102-805270-BP	0.5
102-822855-BP	0.5

<i>catalog_id</i>	$P(\text{catalog\_id})$
1039U	0.1
1082	0.1

<i>keyword</i>	$P(\text{keyword})$
Wapuzange	0.08
Ratschenkabelschere	0.06



<i>manufacturer</i>	<i>keyword</i>	$P(\text{manufacturer}   \text{keyword})$
Unbekannt	Ratschenkabelschere	0.1
Primum	Ratschenkabelschere	0.4

<i>set_id</i>	<i>catalog_id</i>	<i>manufacturer</i>	<i>keyword</i>	<i>ek_amount</i>	$P(\text{ek\_amount}   \text{set\_id}, \text{catalog\_id}, \text{manufacturer}, \text{keyword})$
102-822855-BP	4445	Unbekannt	Ratschenkabelschere	152.63	0.3
102-822855-BP	986	Primum	Ratschenkabelschere	141.69	0.4

$$P(\text{set\_id}, \dots, \text{ek\_amount}) = P(\text{set\_id}) \cdots P(\text{ek\_amount} | \text{set\_id}, \dots, \text{manufacturer})$$

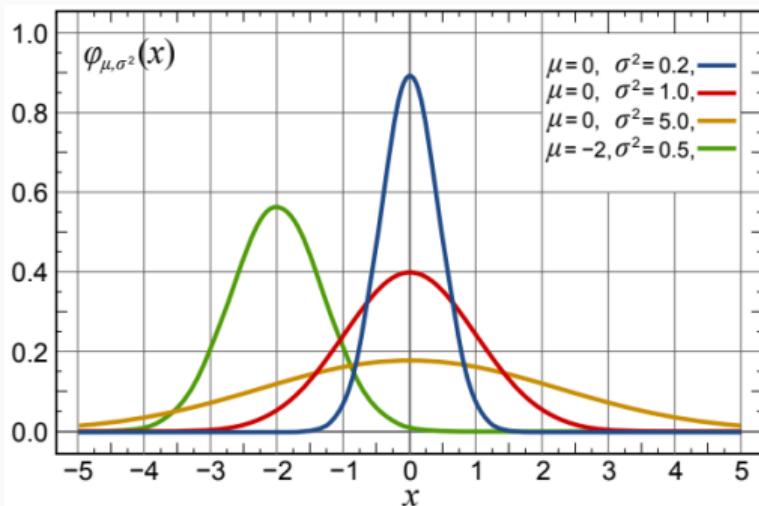
- Directed acyclic graphical model
- Represents a set of **variables** and their **conditional dependencies**

## Domain knowledge:

- “Similar products have similar prices.”
- “Prices might be normally distributed.”

## Normal Distribution ( $\mu \in \mathbb{R}, \sigma > 0$ )

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$



## Domain knowledge:

- “Similar products have similar prices.”
- “Prices might be normally distributed.”

## Normal Distribution ( $\mu \in \mathbb{R}, \sigma > 0$ )

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

## Conditional probability distribution:

$$P(ek\_amount | keyword = "Wapuzange", manufacturer = "Knipex") = \mathcal{N}(ek\_amount | \mu, \sigma^2),$$

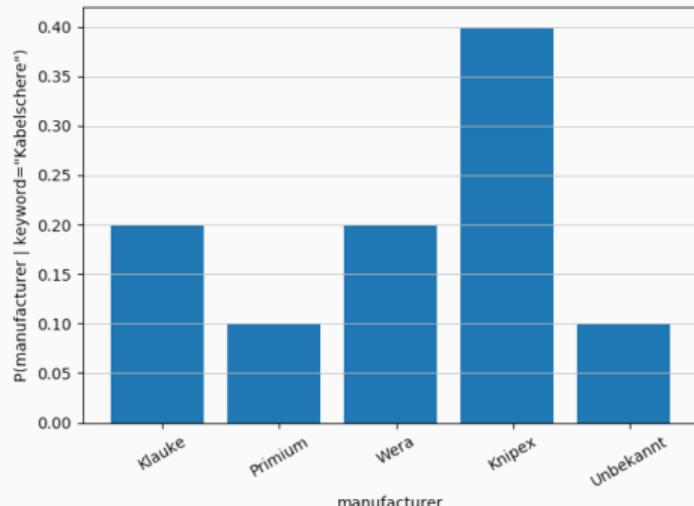
where

$\mu$  ... Mean price of “Knipex Wapuzange” articles

$\sigma^2$  ... Price variance of “Knipex Wapuzange” articles

## Domain knowledge:

- “Different manufacturers make different products.”



## Categorical Distribution

$$P(X = x_i) = p_i, \quad i \in \{1, \dots, k\}$$

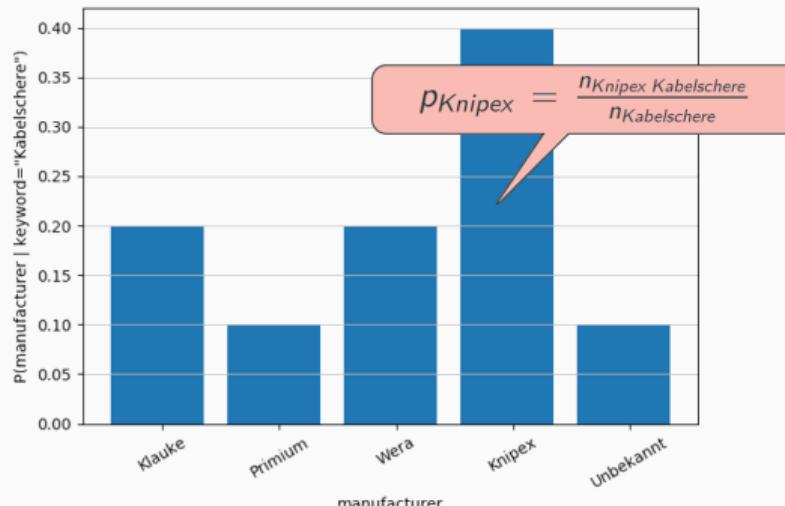
$k > 0$  number of events

$p_1, \dots, p_k$  event probabilities

$$(p_i > 0, \sum p_i = 1)$$

## Domain knowledge:

- “Difference manufacturers make different products.”



## Categorical Distribution

$$P(X = x_i) = p_i, i \in \{1, \dots, k\}$$

$k > 0$  number of events  
 $p_1, \dots, p_k$  event probabilities  
 $(p_i > 0, \sum p_i = 1)$

## Domain knowledge:

- “Different manufacturers make different products.”

## Conditional probability distribution:

$$P(\text{manufacturer} \mid \text{keyword} = \text{"Kabelschere"}) = \text{categorical}(\mathbf{p}),$$

where

$\mathbf{p} = \{p_1, \dots, p_k\}$  ... Probabilities of manufacturers given  $\text{keyword} = \text{"Kabelschere"}$   
⇒ Calculated from occurrence statistics

## Probabilistic Programming:

- Paradigm for specifying probabilistic models



## Our implementation:

- A sampler script to resemble the Bayesian network
- Simulating the data generation process

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

Outlier (?)

```

@gen function mercateo_data_model(observation)
    # No prior knowledge on root nodes => Assume uniform distribution
    set_id = @trace(choose_uniformly(observation.set_id), :set_id)
    catalog_id = @trace(choose_uniformly(observation.catalog_id), :catalog_id)
    keyword = @trace(choose_uniformly(observation.keyword), :keyword)

    # Given keyword, sample manufacturer based on occurrence statistics
    manufacturer = @trace(categorical(observation.manufacturer[keyword]), :manufacturer)

    # Given sampled values above, sample ek_amount based on observed mean and std
    ek_mu = mean(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
    ek_sigma = std(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
    ek_amount = @trace(normal(ek_mu, ek_sigma), :ek_amount)

    return set_id, catalog_id, keyword, manufacturer, ek_amount
end

```

article_id	log( $P(\text{set\_id}, \dots, \text{ek\_amount})$ )
9739011201	-32.7
1528050	-19.7
158-5451	-19.2
201096-YY	-30.2
K-95 31 250	-23.1
0184472	-27.2



set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

Outlier (?)

```

@gen function mercateo_data_model(observation)
    # No prior knowledge on root nodes => Assume uniform distribution
    set_id = @trace(choose_uniformly(observation.set_id), :set_id)
    catalog_id = @trace(choose_uniformly(observation.catalog_id), :catalog_id)
    keyword = @trace(choose_uniformly(observation.keyword), :keyword)

    # Given keyword, sample manufacturer based on occurrence statistics
    manufacturer = @trace(categorical(observation.manufacturer[keyword]), :manufacturer)

    # Given sampled values above, sample ek_amount based on observed mean and std
    ek_mu = mean(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
    ek_sigma = std(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
    ek_amount = @trace(normal(ek_mu, ek_sigma), :ek_amount)

    return set_id, catalog_id, keyword, manufacturer, ek_amount
end

```

article_id	log( $P(\text{set\_id}, \dots, \text{ek\_amount})$ )
9739011201	-32.7
1528050	-19.7
158-5451	-19.2
201096-YY	-30.2
K-95 31 250	-23.1
0184472	-27.2

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

Conditional distribution:

$$P(\text{ek\_amount} \mid \text{keyword} = \text{"Ratschen-Kabelschere"}) = \mathcal{N}(\text{ek\_amount} \mid \mu, \sigma^2)$$



Too many distinct values!

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount	ean
102-805270-BP	1039U	9739011201	Knipex	Eckschwedenzange	C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

## Conditional distribution:

$$P(\text{ek\_amount} \mid \text{keyword} = \text{"Ratschen-Kabelschere"}) = \mathcal{N}(\text{ek\_amount} \mid \mu, \sigma^2)$$

⇒ Not enough evidence to justify  $\mu$  and  $\sigma$ !

set_id	catalog_id	Evaluate $\mu$ and $\sigma$ from similar keywords			unit	ek_amount	ean
102-805270-BP	1039U	9739011	Knickerschweidenzange		C62	55.7	4003773022022
102-805270-BP	1082	1528050	Knipex	Wapuzange	C62	14.82	4003773022022
102-805270-BP	108EL	158-5451	Knipex	Eck-Rohrzange	C62	18.33	
102-822855-BP	4445	201096-YY	Unbekannt	Ratschenkabelschere	C62	152.63	4003773043935
102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0	
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69	4003773043935

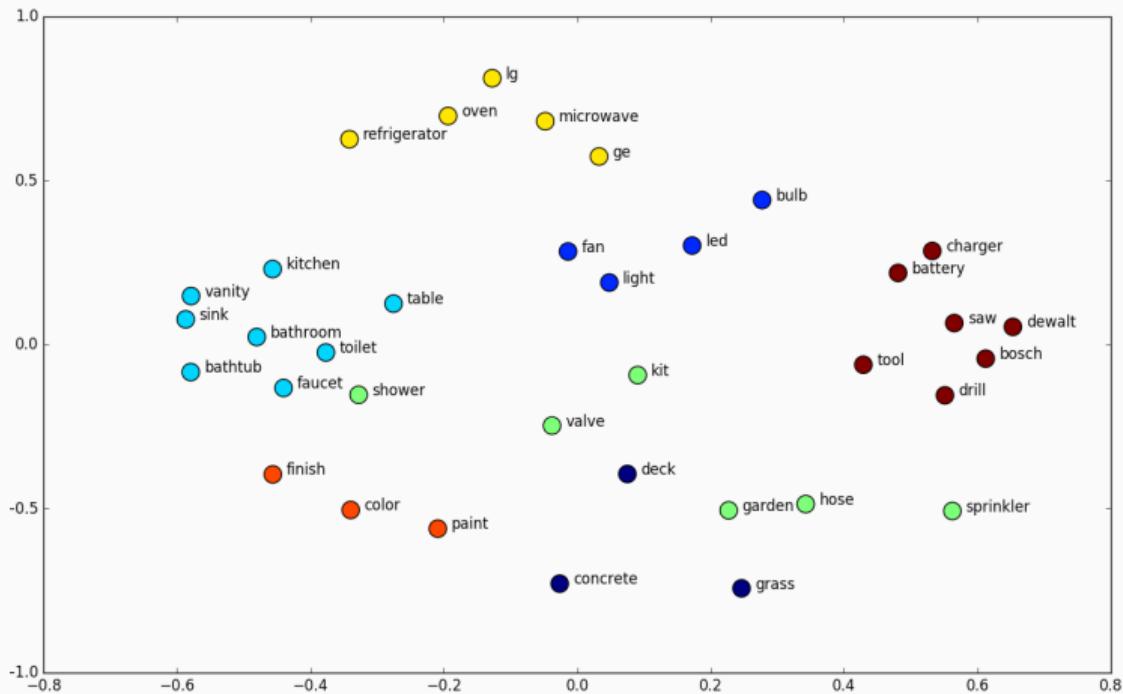
## Conditional distribution:

$$P(\text{ek\_amount} \mid \text{keyword} = \text{"Ratschen-Kabelschere"}) = \mathcal{N}(\text{ek\_amount} \mid \mu, \sigma^2)$$

⇒ Not enough evidence to justify  $\mu$  and  $\sigma$ !

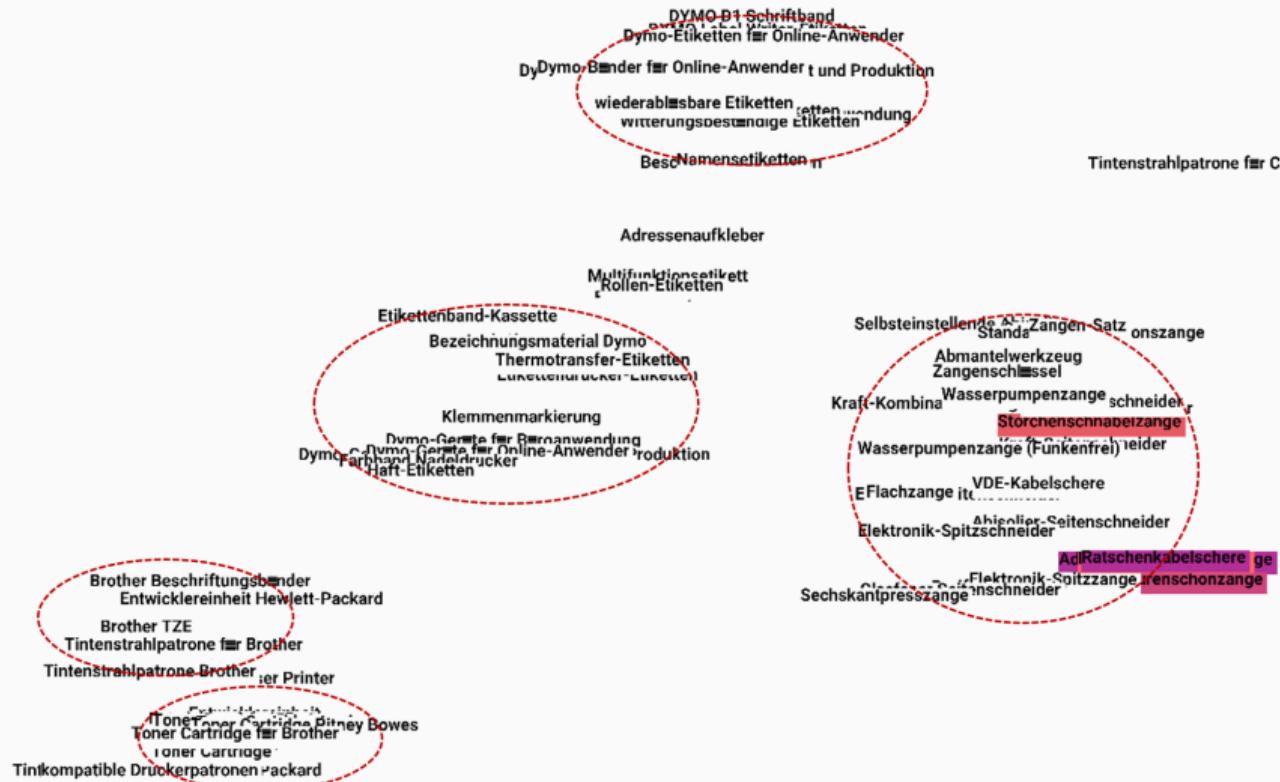
# What are Embeddings?

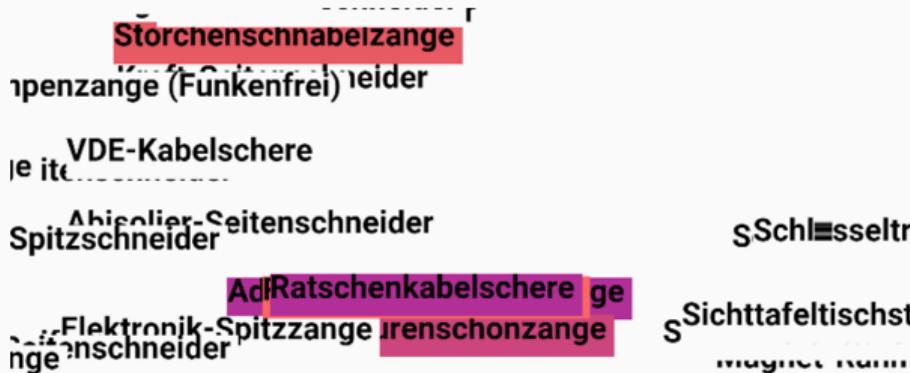
- In NLP: High-dimensional vector representations for words
- Words that share common contexts are close to each other in the vector space





# Neighbors in High-dimensional Space



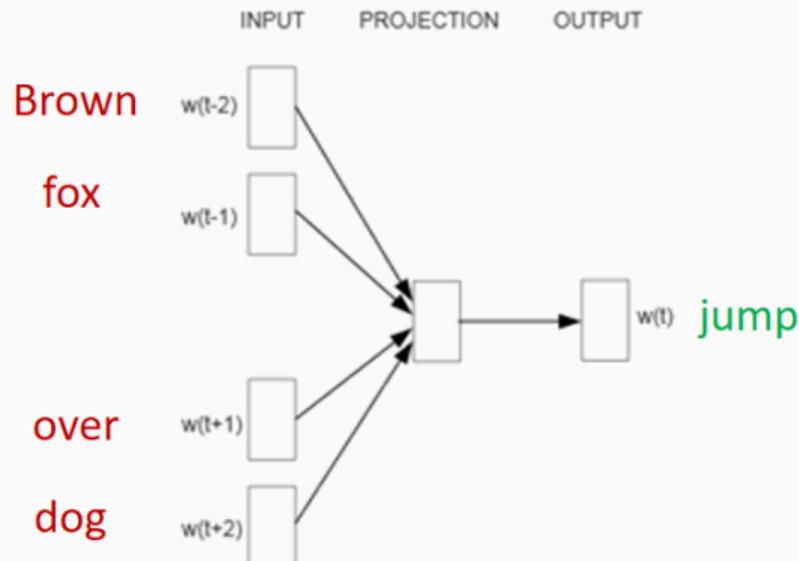


Nearest points in the original space:	
Kabelschneider	0.072
Ratschenkabelschere	0.110
Aderendhülsen-Crimpzange	0.110
Armaturenschonzange	0.198
Storchenschnabelzange	0.282

Conditional distribution:

$$P(\text{ek\_amount} | \text{keyword} = \text{"Ratschen-Kabelschere"}) \\ = P(\text{ek\_amount} | \text{keyword} \in \{\text{"Ratschen-Kabelschere"}, \text{"Kabelschneider"}, \dots\})$$

- Derived from *Word2vec*<sup>2</sup> models
- Given “**context**”, predict “**target**”



<sup>2</sup>Mikolov, Tomas, et al. “Efficient estimation of word representations in vector space.”

... The quick brown fox jumps over the lazy dog. ...

102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0
---------------	-----	-------------	--------	----------------------	-----	-------

## Similarities:

- Sequences with contextual information
- Limited knowledge source for high-dimensional mapping

... The quick brown fox jumps over the lazy dog. ...

102-822855-BP	381	K-95 31 250	Knipex	Ratschen-Kabelschere	C62	118.0
---------------	-----	-------------	--------	----------------------	-----	-------

## Similarities:

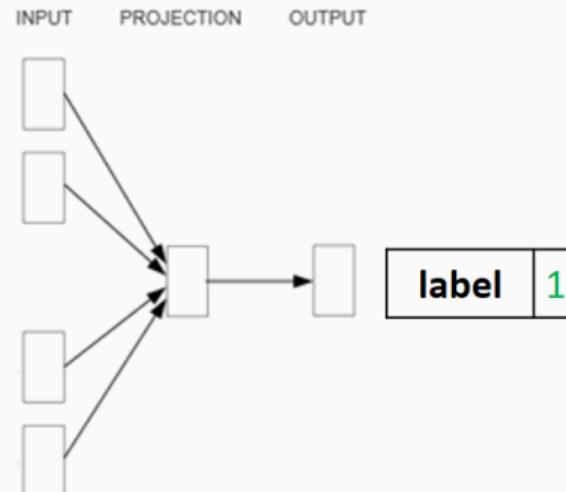
- Sequences with contextual information
- Limited knowledge source for high-dimensional mapping

## Differences:

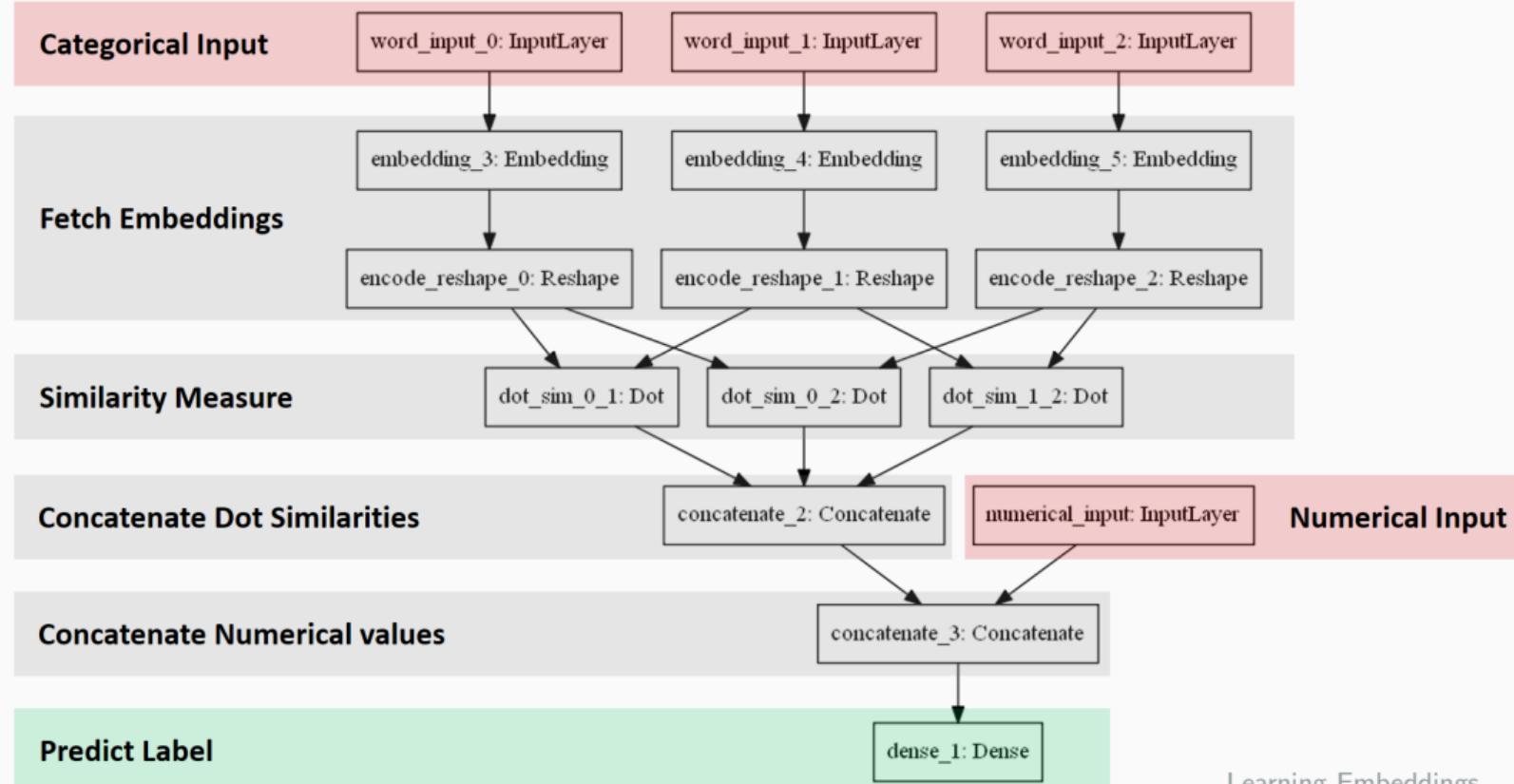
- Variable vs. fixed sequence length
- Single vs. multiple vocabularies
- Numerical values in tabular data

- Derived from *Word2vec*<sup>3</sup> models
- Given “context”, predict “target”

<b>set_id</b>	102-822855-BP
<b>catalog_id</b>	381
...	...
<b>keyword</b>	Ratschen-Kabelschere
<b>ek_amount</b>	118.0



<sup>2</sup>Mikolov, Tomas, et al. “Efficient estimation of word representations in vector space.”



## To investigate:

- Effectiveness of probabilistic model
- Benefits from embeddings regarding anomaly detection performance

No ground  
truth  
information  
available

## To investigate:

- Effectiveness of probabilistic model
- Benefits from embeddings regarding anomaly detection performance

No ground  
truth  
information  
available

## To investigate:

- Effectiveness of probabilistic model
- Benefits from embeddings regarding anomaly detection performance

## Procedure:

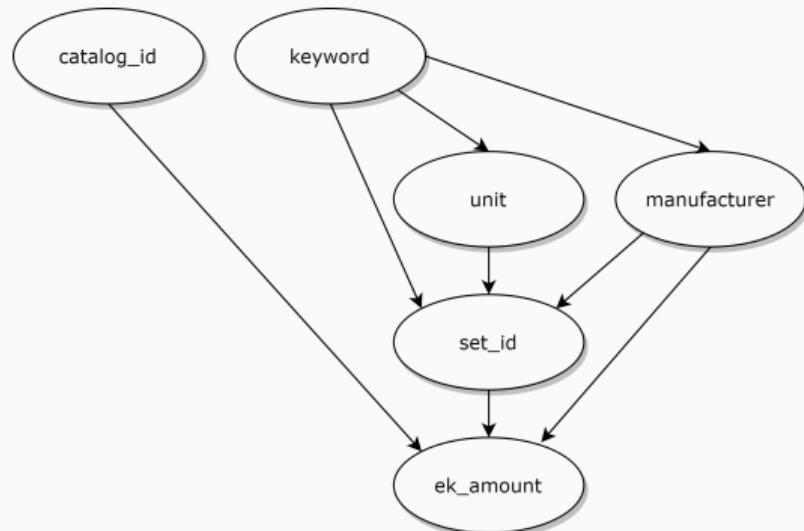
1. Consider observed data as “normal”
2. Manually generate “anomalous” tuples
3. Build probabilistic model and train embeddings using combined data
4. Evaluate confidence scores
5. Check if manual anomalies have low confidence scores

## Mercateo article database

- Top 200 sets of articles
- 9049 data tuples as observation
- 452 manually generated anomalies (5% of observation)

Attribute	Attribute type	Cardinality
<i>catalog_id</i>	Categorical	210
<i>unit</i>	Categorical	10
<i>keyword</i>	Categorical	167
<i>manufacturer</i>	Categorical	35
<i>set_id</i>	Categorical	200
<i>ek_amount</i>	Numerical	n.a.

## Bayesian network:



## Precision

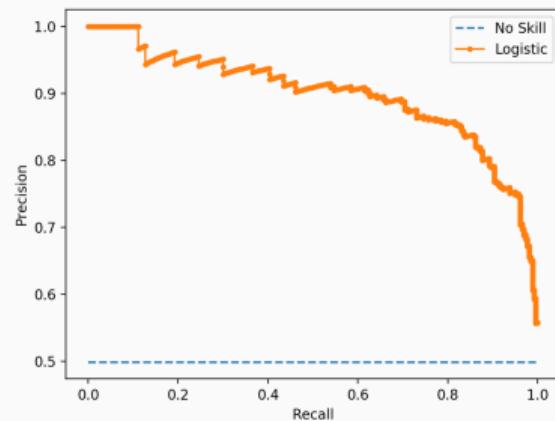
$$Precision = \frac{TP}{TP + FP}$$

## Recall

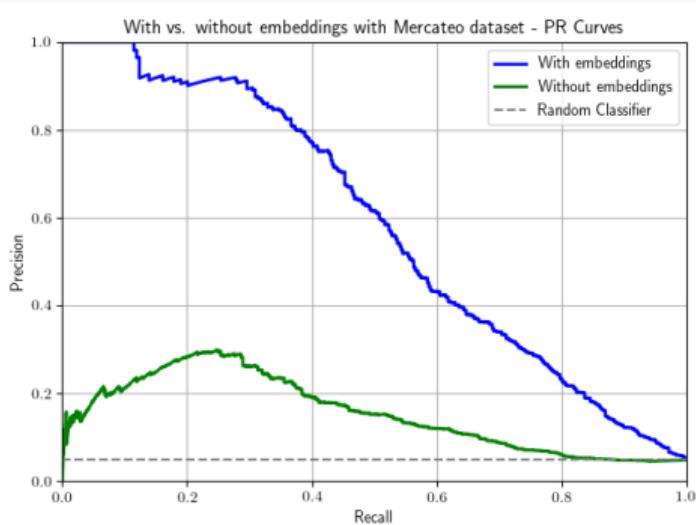
$$Recall = \frac{TP}{TP + FN}$$

## Precision-recall curve:

- Describes how good a model predicts positive class (anomalies)
- Goal: Keep precision high while recall increases

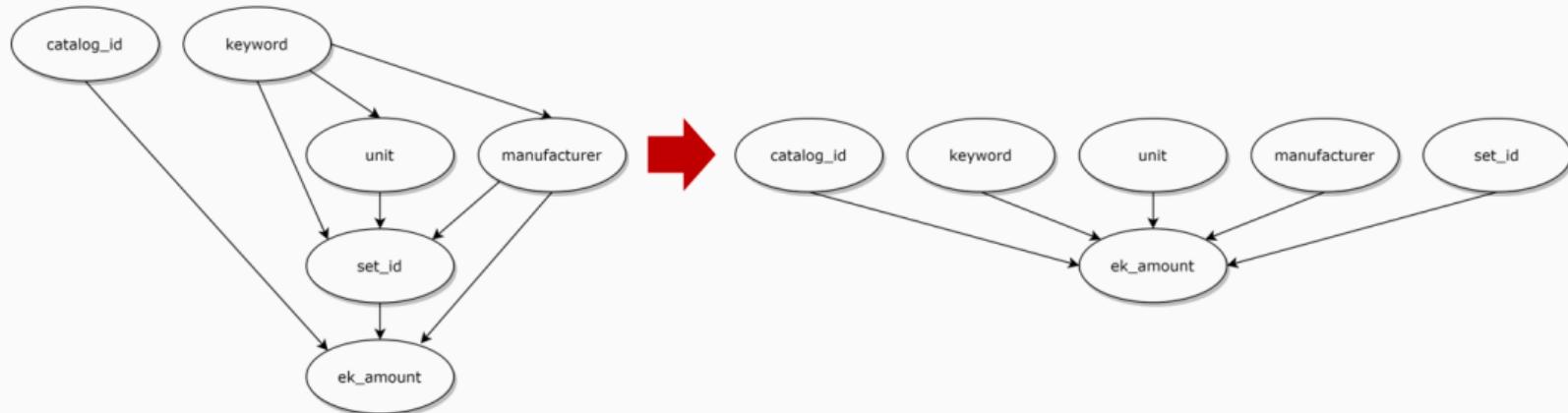


## Results - With vs. Without Embeddings

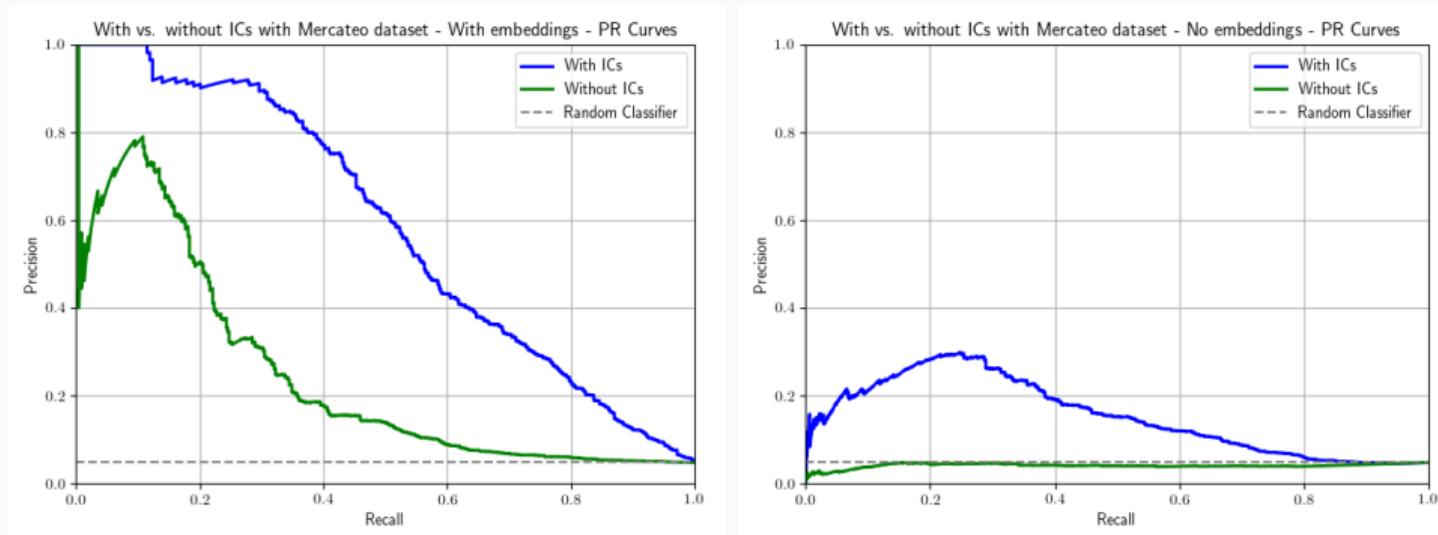


	<b>With embeddings</b>	<b>Without embeddings</b>
Recall@ $P = 452$	0.54	0.29

# Results - Minimizing External Knowledge



# Results - Minimizing External Knowledge



	<b>With embeddings</b>	<b>Without embeddings</b>
Recall@P = 452	0.54	0.29
Recall@P = 452, minimal EK	0.30	0.02

Embeddings are helpful considering

- Improvement in anomaly detection under high-cardinality
- Compensation for limited external knowledge

**Anomaly detection is feasible with probabilistic models!**

## Future steps:

- Further examination of the neural network model
- Experiments with larger datasets
- Automated repair suggestion

**Thank you! (Questions?)**

```
1 @gen function mercateo_data_model(observation)
2     # No prior knowledge on root nodes => Assume uniform distribution
3     set_id = @trace(choose_uniformly(observation.set_id), :set_id)
4     catalog_id = @trace(choose_uniformly(observation.catalog_id), :catalog_id)
5     keyword = @trace(choose_uniformly(observation.keyword), :keyword)
6
7     # Given keyword, sample manufacturer based on occurrence statistics
8     manufacturer = @trace(categorical(observation.manufacturer[keyword]), :manufacturer)
9
10    # Given sampled values above, sample ek_amount based on observed mean and std
11    ek_mu = mean(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
12    ek_sigma = std(observation.ek_amount[set_id, catalog_id, keyword, manufacturer])
13    ek_amount = @trace(normal(ek_mu, ek_sigma), :ek_amount)
14
15    return set_id, catalog_id, keyword, manufacturer, ek_amount
16 end
```

```
1 # Original keyword
2 keyword = @trace(choose_uniformly(observation.keyword), :keyword)
3
4 # Find neighboring keywords
5 keyword_neighborhood = find_neighbors(keyword, keyword_embeddings)
6
7 #  $P(\text{manufacturer} | \text{keyword} \in \text{keyword\_neighborhood})$ 
8 manufacturer = @trace(categorical(observation.manufacturer[keyword_neighborhood]), :manufacturer)
9
10
11 # Find neighboring manufacturers
12 manufacturer_neighborhood = find_neighbors(manufacturer, manufacturer_embeddings)
13
14 #  $P(\text{ek\_amount} | \text{keyword} \in \text{keyword\_neighborhood}, \text{manufacturer} \in \text{manufacturer\_neighborhood})$ 
15 ek_mu = mean(observation.ek_amount[keyword_neighborhood, manufacturer_neighborhood])
16 ek_sigma = std(observation.ek_amount[keyword_neighborhood, manufacturer_neighborhood])
17 ek_amount = @trace(normal(ek_mu, ek_sigma), :ek_amount)
```

## Examples from “False Positives”

set_id	catalog_id	article_id	manufacturer	keyword	unit	ek_amount
102-805270-BP	1039U	9739011201	Knipex	Eckschwendenzange	C62	55.7
102-822855-BP	986	0184472	Primum	Ratschenkabelschere	C62	141.69
1029-1017005	1039U	9859711444	Brother	Etiketten-Papier	PK	20.45
102-823785-BP	1039U	9739011068	Knipex	Stillson-Rohrzange	C62	42.58
102-801909-BP	1039U	9739011049	Knipex	Abmantelwerkzeug	C62	80.7