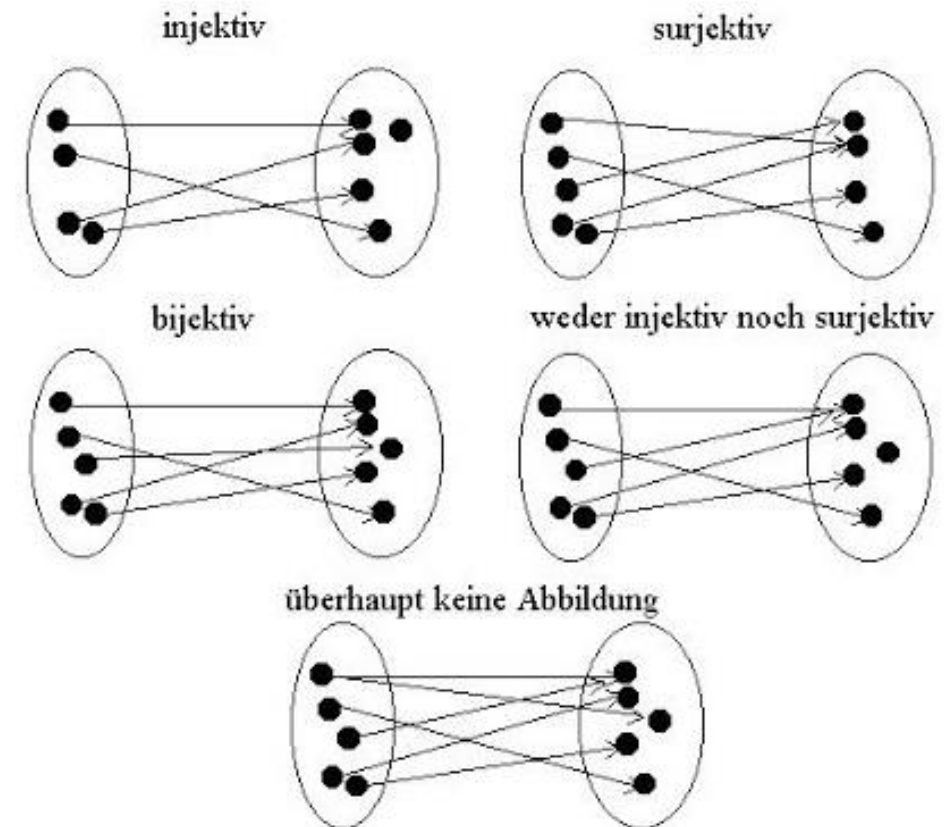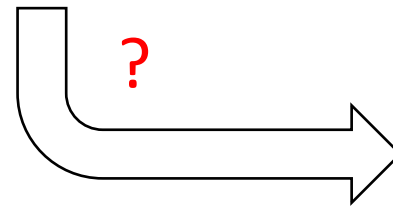# Remark: Text Preprocessing

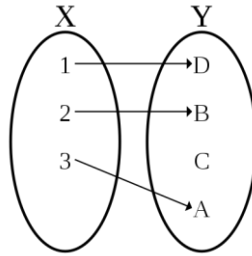F(w) := Text Preprocessing

- Tokenization
- Character Set
- Punctuation
- ....

# Injective, surjective, bijective - What??

$f: X \rightarrow Y$ is called          $\Leftrightarrow$

?                      injective

# Injective, surjective, bijective - What??



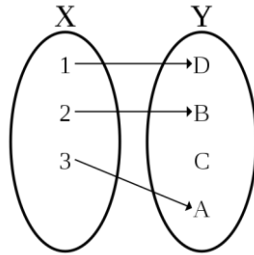$f: X \rightarrow Y$ is called                    $\Leftrightarrow$

injective                    $x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$

?                    surjective

Technische **Hochschule**
Ingolstadt

# Injective, surjective, bijective - What??

$f : X \rightarrow Y$ is called $\qquad\qquad \Leftrightarrow$

injective $\qquad\qquad x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$

surjective $\qquad\qquad \forall y \in Y \; \exists x \in X : f(x) = y$

? $\qquad\qquad$ bijective

# Injective, surjective, bijective - What??



$f : X \to Y$ is called $\qquad\qquad \Leftrightarrow$

injective $\qquad\qquad x_1 \neq x_2 \Rightarrow f(x_1) \neq f(x_2)$

surjective $\qquad\qquad \forall y \in Y \; \exists x \in X : f(x) = y$

bijective $\qquad\qquad$ f is injective and surjective

# Remark: Text Preprocessing

F(w) := Text Preprocessing

- Tokenization
- Character Set
- Punctuation
- ....

?

# Remark: Text Preprocessing

F(w) := Text Preprocessing

- Tokenization
- Character Set
- Punctuation
- ....



injektiv

surjektiv

bijektiv

weder injektiv noch surjektiv

überhaupt keine Abbildung

Desire/Goal…

Reality…

# Remark: Text Preprocessing

F($\underline{w}$) := Text Preprocessing

- Tokenization
- Character Set
- Punctuation
- ....



Desire/Goal…

Reality in big projects …

Reality…

How to weigh the *association between two words*?

How much more do the two words co-occur in our corpus

than we would have a priori expected them to appear by chance?

**Intuition**: „How much more do the two words co-occur in our corpus than we would have a priori expected them to appear by chance?"

**Definition**

The PMI between target word $w$ and context word $c$ is defined as

$$PMI(w,c) := \log_2 \frac{P(w,c)}{P(w)P(c)}.$$

Further Reading: Chapter 6.6, https://web.stanford.edu/~jurafsky/slp3/

$$PMI(w,c) := \log_2 \frac{P(w,c)}{P(w)P(c)}$$

```
PMI >> 0        := co-occurrence of wᵢ and wⱼ
PMI near zero  := wᵢ and wⱼ are unlikely as word pair in the corpus
```

# Pointwise Mutual Information

## Examples

$$PMI(w, c) := \log_2 \frac{P(w, c)}{P(w)P(c)}$$

$$PMI(\text{santa}, \text{fe}) = 13.74529$$

$$PMI(\text{hong}, \text{kong}) = 12.54758$$

$$PMI(\text{las}, \text{vegas}) = 12.53122$$

$$PMI(\text{los}, \text{angeles}) = 11.16868$$

$$PMI(\text{tele}, \text{communications}) = 8.13129$$

$$PMI(\text{multimillion}, \text{dollar}) = 8.03381$$

$$PMI(\text{U.}, \text{S.}) = 7.21532.$$

# Pointwise Mutual Information

*Examples*

$$PMI(w, c) := \log_2 \frac{P(w, c)}{P(w)P(c)}$$

$$PMI(\text{santa}, \text{fe}) = 13.74529$$

$$PMI(\text{hong}, \text{kong}) = 12.54758$$

$$PMI(\text{las}, \text{vegas}) = 12.53122$$

$$PMI(\text{los}, \text{angeles}) = 11.16868$$

$$PMI(\text{tele}, \text{communications}) = 8.13129$$

$$PMI(\text{multimillion}, \text{dollar}) = 8.03381$$

$$PMI(\text{U.}, \text{S.}) = 7.21532.$$

$$PMI(\text{accounted}, \text{for}) = 5.98041$$

$$PMI(\text{according}, \text{to}) = 5.93769$$

$$PMI(\text{intends}, \text{to}) = 5.9195$$

$$PMI(\text{able}, \text{to}) = 5.91757$$

# Pointwise Mutual Information

*Examples*

$$PMI(w, c) := \log_2 \frac{P(w, c)}{P(w)P(c)}$$

$$PMI(\text{manufacturer}, \text{holds}) = 0.00058$$

$$PMI(\text{stocks}, \text{shows}) = 0.00013$$

$$PMI(\text{legal}, \text{confrontation}) = -0.00013$$

$$PMI(\text{pennsylvania}, \text{cleveland}) = -0.00061$$

$$PMI(w, c) := \log_2 \frac{P(w, c)}{P(w)P(c)}$$

PMI(manufacturer, holds) = 0.00058

PMI(stocks, shows) = 0.00013

PMI(legal, confrontation) = −0.00013

PMI(pennsylvania, cleveland) = −0.00061

PMI(not, international) = −4.50398

PMI(two, interests) = −5.00358

PMI(the, said) = −8.95582

# Pointwise Mutual Information

*Interpretation*

$$PMI(w, c) := \log_2 \frac{P(w, c)}{P(w)P(c)}$$

```
PMI >> 0        := co-occurrence of wᵢ and wⱼ
PMI near zero   := wᵢ and wⱼ are unlikely as word pair in the corpus
```

$$PMI(w, c) \stackrel{?}{=} PMI(c, w)$$

# Mutual Information

"It quantifies the amount of information obtained about one random variable through observing the other random variable."

# What was entropy?

$$H(X) = \boxed{\text{Who remembers it?}}$$

# What was entropy?

$$H(X) = -\sum_{w \in X} p(w) \log_2 p(w)$$

# What was conditional entropy?

$$H(Y|X) = \sum_{w \in X} p(w) H(Y|X = w)$$

$$= \sum_{w \in X} p(w) [-\sum_{v \in Y} p(v|w) \log_2 p(v|w)]$$

Sprach- und Textverstehen 2020
Prof. Dr. Munir Georges

Technische Hochschule
Ingolstadt

# Mutual Information

In Deutsch: Transinformation

$H(y|x)$

$H(x|y)$

$H(x)$

$H(y)$

$I(x,y)$

Channel

Input $x$

Output $y = x + noise$

H(X)

H(Y)

H(X|Y)  I(X;Y)  H(Y|X)

H(X,Y)

# Mutual Information

In Deutsch: Transinformation



$$I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

# Mutual Information

In Deutsch: Transinformation

$$I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

# Mutual Information

### In Deutsch: Transinformation

$$H(Y|X) = \sum_{w \in X} p(w) H(Y|X = w)$$

$$= \sum_{w \in X} p(w) [- \sum_{v \in Y} p(v|w) \log_2 p(v|w)]$$

$$H(X) = - \sum_{w \in X} p(w) \log_2 p(w)$$

$$I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

# Mutual Information

### In Deutsch: Transinformation

$$I(w_I, w_J) := \sum_{i,j} P(w_i, w_j) \log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}$$

$$H(Y|X) = \sum_{w \in X} p(w)H(Y|X = w)$$

$$= \sum_{w \in X} p(w)[-\sum_{v \in Y} p(v|w) \log_2 p(v|w)]$$

$$H(X) = -\sum_{w \in X} p(w) \log_2 p(w)$$

$$I(X;Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

# Remarks

Pointwise Mutual Information

$$I(w_I, w_J) := \sum_{i,j} P(w_i, w_j) \log \frac{p(w_i, w_j)}{p(w_i) p(w_j)}$$

- "In many applications, one wants to **maximize mutual information** (thus increasing dependencies), which is often equivalent to **minimizing conditional entropy**."



Further Reading: https://en.wikipedia.org/wiki/Mutual_information

# Remarks

Pointwise Mutual Information

$$I(w_I, w_J) := \sum_{i,j} P(w_i, w_j) \log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}$$

- "In many applications, one wants to **maximize mutual information** (thus increasing dependencies), which is often equivalent to **minimizing conditional entropy**."



Question: Relation to Language Models or ML in general?
- LMs: minimize perplexity ⇔ maximize probability of test set ⇔ minimize cross entropy (-> Jurafsky, Section 3.8)
- ML: goal is often: minimize cross entropy loss

Further Reading: https://en.wikipedia.org/wiki/Mutual_information

Sprach- und Textverstehen 2020
Prof. Dr. Munir Georges

Technische Hochschule
Ingolstadt

# Remarks

Pointwise Mutual Information

$$I(w_I, w_J) := \sum_{i,j} P(w_i, w_j) \log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}$$

- "In many applications, one wants to **maximize mutual information** (thus increasing dependencies), which is often equivalent to **minimizing conditional entropy**."

- "mutual information between phrases and contexts is used as a feature for k-means clustering to discover semantic clusters (concepts)"

$H(x)$  $H(x|y)$  $H(y|x)$  $H(y)$

$I(x,y)$

Channel

Input $x$    Output $y = x + noise$

Further Reading: https://en.wikipedia.org/wiki/Mutual_information

Sprach- und Textverstehen 2020
Prof. Dr. Munir Georges

Technische Hochschule
Ingolstadt

# Remarks

Pointwise Mutual Information

$$\mathrm{I}(w_I, w_J) := \sum_{i,j} P(w_i, w_j) \log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}$$

- "In many applications, one wants to **maximize mutual information** (thus increasing dependencies), which is often equivalent to **minimizing conditional entropy**."
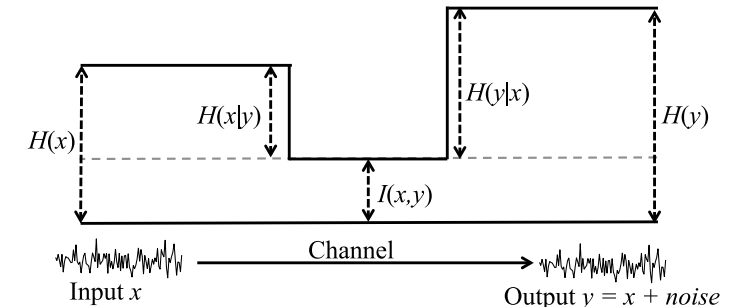
- "mutual information between phrases and contexts is used as a feature for k-means clustering to discover semantic clusters (concepts)"

- Can be used for feature selection in Machine Learning:
  minimum Redundancy Maximum Relevance (mRMR) Algorithm

$H(x)$  $H(x|y)$  $H(y|x)$  $H(y)$

$I(x,y)$

Channel

Input $x$  Output $y = x + noise$

Further Reading: https://en.wikipedia.org/wiki/Mutual_information    Further Watching: https://youtu.be/YKDZHPJ-pQ0

# Text correction

# Text correction

| Input | To check | Correction |
|---|---|---|
| I don't know **whether** I want to | **whether** | whether |
| The **weather** is pretty bad today | **weather** | weather |
| **Whether** he likes me or not I can't say | **whether** | whether |
| I like sunny **weather** | **weather** | weather |

Any Idea?

https://en.wikipedia.org/wiki/Levenshtein_distance

# Text correction

| Input | To check | Correction |
|---|---|---|
| I don't know **whether** I want to | **whether** | whether |
| The **weather** is pretty bad today | **weather** | weather |
| **Whether** he likes me or not I can't say | **whether** | whether |
| I like sunny **weather** | **weather** | weather |

Levenshtein distance("whether", "weather")

…

# Text correction

| Input | To check | Correction |
|---|---|---|
| I don't know **whether** I want to | **whether** | whether |
| The **weather** is pretty bad today | **weather** | weather |
| **Whether** he likes me or not I can't say | **whether** | whether |
| I like sunny **weather** | **weather** | weather |

Idea?

https://www.cl.uni-heidelberg.de/courses/ws19/ecl/material/ngram_handout.pdf

# Text correction

| Input | To check | Correction |
|---|---|---|
| I don't know **whether** I want to | **whether** | whether |
| The **weather** is pretty bad today | **weather** | weather |
| **Whether** he likes me or not I can't say | **whether** | whether |
| I like sunny **weather** | **weather** | weather |

P(know whether) vs. P(know weather)

P(whether I) vs. P(weather I)

…

https://www.cl.uni-heidelberg.de/courses/ws19/ecl/material/ngram_handout.pdf

Technische Hochschule
Ingolstadt

# Code-Switch Detection

**„Wenn Sie dort sind, please do not talk about our product. Ist das klar?"**

https://www.csl.uni-bremen.de/cms/images/documents/publications/BA_SebastianLeidig.pdf

# Code-Switch Detection

„**Wenn Sie dort sind**, please do not talk about our product. **Ist das klar?**"

↓           ↓           ↓

**German**           English           **German**

https://www.csl.uni-bremen.de/cms/images/documents/publications/BA_SebastianLeidig.pdf

# Anglicism Detection

„Mach ein **mind map**"  → English

vs.

„Mach eine **Gedächtniskarte**"  → German

„ **Cancel** mal die Bestellung"  → English

vs.

„**Storniere** mal die Bestellung"  → German (from Latin)

vs.

„**Mach** mal die Bestellung **rückgängig**"  → German

https://www.csl.uni-bremen.de/cms/images/documents/publications/BA_SebastianLeidig.pdf

# Problems

**Pseudo-Anglicism**

German: Handy vs. English: Mobile Phone

# Problems

**Pseudo-Anglicism**

German: Handy vs. English: Mobile Phone

German: Beamer vs. English: Projector

# Problems

**Pseudo-Anglicism**

German: Handy vs. English: Mobile Phone

German: Beamer vs. English: Projector

**Hybrid word**

Aquaphobia: Latin "aqua" (water), Greek "phobia"(fear)

https://www.csl.uni-bremen.de/cms/images/documents/publications/BA_SebastianLeidig.pdf

# Problems

**Pseudo-Anglicism**

German: Handy vs. English: Mobile Phone

German: Beamer vs. English: Projector

**Hybrid word**

Aquaphobia: Latin "aqua" (water), Greek "phobia"(fear)

Schadsoftware: German "Schaden" (damage)

# Problems

**Pseudo-Anglicism**

German: Handy vs. English: Mobile Phone

German: Beamer vs. English: Projector

**Hybrid word**

Aquaphobia: Latin "aqua" (water), Greek "phobia"(fear)

Schadsoftware: German "*Schad*en" (damage)

⤷ English: malware

…

https://www.csl.uni-bremen.de/cms/images/documents/publications/BA_SebastianLeidig.pdf

Technische Hochschule Ingolstadt

# Naive Bayes Classifiers

## Naive Bayes for text classification

# What does a NB classifier model look like?

$$\hat{c} = \text{argmax}_{c \in C} P(c|d)$$
$$= \text{argmax}_{c \in C} P(d|c)P(c)$$

# What does a NB classifier model look like?

$$\hat{c} = \text{argmax}_{c \in C} P(c|d)$$
$$= \text{argmax}_{c \in C} P(d|c) P(c)$$

Compute most probable class, given some document d
by

       choosing the class which has the highest product of two
       probabilities

# What does a NB classifier model look like?

$$\hat{c} = \text{argmax}_{c \in C} P(c|d)$$
$$= \text{argmax}_{c \in C} P(d|c)P(c)$$

Compute most probable class,  given some document d
by
choosing the class which has the highest product  of two
probabilities:
likelihood of document and class prior probability

# How to represent a document?

Remember: $\hat{c} = \text{argmax}_{c \in C} P(d|c)P(c)$

document $d$ ⟷ set of features $f_1, \dots, f_n$

$$\hat{c} = \text{argmax}_{c \in C} P(f_1, \dots, f_n|c)P(c)$$

Technische Hochschule Ingolstadt

# Two simplifying assumptions

# Two simplifying assumptions

1. A word's position doesn't matter

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)

2. Features $\{f_1, \dots, f_n\} = F$ only encode word identity, not its position

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)

2. Features $\{f_1, \ldots, f_n\} = F$ only encode word identity, not its position

2nd assumption ⇔ **Naive Bayes Assumption**

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)

2. Features $\{f_1, \ldots, f_n\} = F$ only encode word identity, not its position

2nd assumption $\Leftrightarrow$ **Naive Bayes Assumption**

$$P(f_1, \ldots, f_n | c) = P(f_1 | c) \cdot P(f_2 | c) \cdot \ldots \cdot P(f_n | c)$$

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)
2. Features $\{f_1, \ldots, f_n\} = F$ only encode word identity, not its position

2nd assumption $\Leftrightarrow$ ***Naive Bayes Assumption***
*= conditional independence assumption*

$$P(f_1, \ldots, f_n | c) = P(f_1 | c) \cdot P(f_2 | c) \cdot \ldots \cdot P(f_n | c)$$

# Two simplifying assumptions

1. A word's position doesn't matter   (Bag-of-Words assumption)
2. Features $\{f_1, \dots, f_n\} = F$ only encode word identity, not its position

2<sup>nd</sup> assumption $\Leftrightarrow$ **Naive Bayes Assumption**
= *conditional independence assumption*

$$P(f_1, \dots, f_n | c) = P(f_1|c) \cdot P(f_2|c) \cdot \dots \cdot P(f_n|c)$$

$$\Rightarrow \hat{c} = \mathrm{argmax}_{c \in C}\ P(c) \cdot \prod_{f \in F} P(f|c)$$

Technische Hochschule
Ingolstadt

# Calculating NB document classifier

Remember: $\hat{c} = \text{argmax}_{c \in C} \, P(c) \cdot \prod_{f \in F} P(\textcolor{green}{f}|c)$

# Calculating NB document classifier

Remember: $\quad\quad\quad \hat{c} = \text{argmax}_{c \in C} \, P(c) \cdot \prod_{f \in F} P(f|c)$

Let $X$ be the set of word positions in test document $d = (w_1, \dots, w_{|X|})$.

# Calculating NB document classifier

Remember: $\qquad \hat{c} = \text{argmax}_{c \in C} \, P(c) \cdot \prod_{f \in F} P(f|c)$

Let $X$ be the set of word positions in test document $d = (w_1, \dots, w_{|X|})$. Then compute

$$\hat{c} = \text{argmax}_{c \in C} \, P(c) \cdot \prod_{i \in X} P(w_i|c)$$

# Calculating NB document classifier

Remember: $\hat{c} = \text{argmax}_{c \in C}\ P(c) \cdot \prod_{f \in F} P(f|c)$

Let $X$ be the set of word positions in test document $d = (w_1, \ldots, w_{|X|})$. Then compute

$$\hat{c} = \text{argmax}_{c \in C}\ P(c) \cdot \prod_{i \in X} P(w_i|c)$$

Question: How to increase speed and avoid underflow?

# Calculating NB document classifier

Remember: $\qquad \hat{c} = \text{argmax}_{c \in C} \; P(c) \cdot \prod_{f \in F} P(f|c)$

Let $X$ be the set of word positions in test document $d = (w_1, \ldots, w_{|X|})$. Then compute

$$\hat{c} = \text{argmax}_{c \in C} \; P(c) \cdot \prod_{i \in X} P(w_i|c)$$

Question: How to increase speed and avoid underflow?

$$\hat{c} = \text{argmax}_{c \in C} \; \log P(c) + \sum_{i \in X} \log P(w_i|c)$$

Technische Hochschule
Ingolstadt

# Calculating NB document classifier

Remember: $\quad \hat{c} = \mathrm{argmax}_{c \in C} \; P(c) \cdot \prod_{f \in F} P(f|c)$

Let $X$ be the set of word positions in test document $d = (w_1, \dots, w_{|X|})$.
Then compute

$$\hat{c} = \mathrm{argmax}_{c \in C} \; P(c) \cdot \prod_{i \in X} P(w_i|c)$$

Question: How to increase speed and avoid underflow?

$$\hat{c} = \mathrm{argmax}_{c \in C} \; \log P(c) + \sum_{i \in X} \log P(w_i|c)$$

$$\hat{P}(w_i|c) = \frac{count(w_i, c)}{\sum_{w \in V} count(w, c)}$$

# Remark on Naive Bayes Classification

- NB classifier is a supervised machine learning method $f: X \to Y$

# Remark on Naive Bayes Classification

- NB classifier is a supervised machine learning method $f\colon X \rightarrow Y$
  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \dots, c_N\}$
  - Goal: return predicted class $\hat{c} \in Y$

# Remark on Naive Bayes Classification

Remember: $\qquad f(d) = \text{argmax}_c \, P(d|c)P(c)$

- NB classifier is a supervised machine learning method $f: X \rightarrow Y$
  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \dots, c_N\}$
  - Goal: return predicted class $\hat{c} \in Y$

# Remark on Naive Bayes Classification

Remember: $\qquad f(d) = \text{argmax}_c \, P(d|c)P(c)$

- NB classifier is a supervised machine learning method $f: X \rightarrow Y$

  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \ldots, c_N\}$

  - Goal: return predicted class $\hat{c} \in Y$

- Given an observation, a NB classifier returns the class most likely to have generated the observation $\Leftrightarrow$ **generative model**

# Remark on Naive Bayes Classification

Remember: $\quad\quad$ $f(d) = \text{argmax}_c \, P(d|c)P(c)$

- NB classifier is a supervised machine learning method $f: X \rightarrow Y$
  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \ldots, c_N\}$
  - Goal: return predicted class $\hat{c} \in Y$

- Given an observation, a NB classifier returns the class most likely to have generated the observation $\Leftrightarrow$ **generative model**
  - first a class is sampled from P(c)

# Remark on Naive Bayes Classification

Remember: $\qquad f(d) = \text{argmax}_c\, P(d|c)P(c)$

- NB classifier is a supervised machine learning method $f: X \rightarrow Y$
  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \ldots, c_N\}$
  - Goal: return predicted class $\hat{c} \in Y$

- Given an observation, a NB classifier returns the class most likely to have generated the observation $\Leftrightarrow$ **generative model**
  - first a class is sampled from P(c)
  - then generate document/words/features… by sampling from P(d|c)

# Remark on Naive Bayes Classification

Remember: $\qquad$ $f(d) = \text{argmax}_c\, P(d|c)P(c)$

- NB classifier is a supervised machine learning method $f: X \to Y$
  - Given: input $x \in X$ and a fixed set of output classes $Y = \{c_1, \ldots, c_N\}$
  - Goal: return predicted class $\hat{c} \in Y$

- Given an observation, a NB classifier returns the class most likely to have generated the observation $\Leftrightarrow$ **generative model**
  - first a class is sampled from P(c)
  - then generate document/words/features… by sampling from P(d|c)

Question: how to generate text in this way?

# Digression
## Sample sentence from a *unigram* LM

Assume we have access to words in $\quad V \cup \{< EOS >\}$

Assume sorted word probabilites: $\quad P(w_1) > P(w_2) > \ ... > P\left(w_{|V|+1}\right)$

# Digression

## Sample sentence from a *unigram* LM

Assume we have access to words in $\qquad V \cup \{< EOS >\}$

Assume sorted word probabilites: $\qquad P(w_1) > P(w_2) > \, ... > P\left(w_{|V|+1}\right)$



**Note:** $\quad 1 = \sum_i P(w_i)$

# Digression

## Sample sentence from a *unigram* LM

Assume we have access to words in $\qquad V \cup \{< EOS >\}$

Assume sorted word probabilites: $\qquad P(w_1) > P(w_2) > \ \ldots > P(w_{|V|+1})$
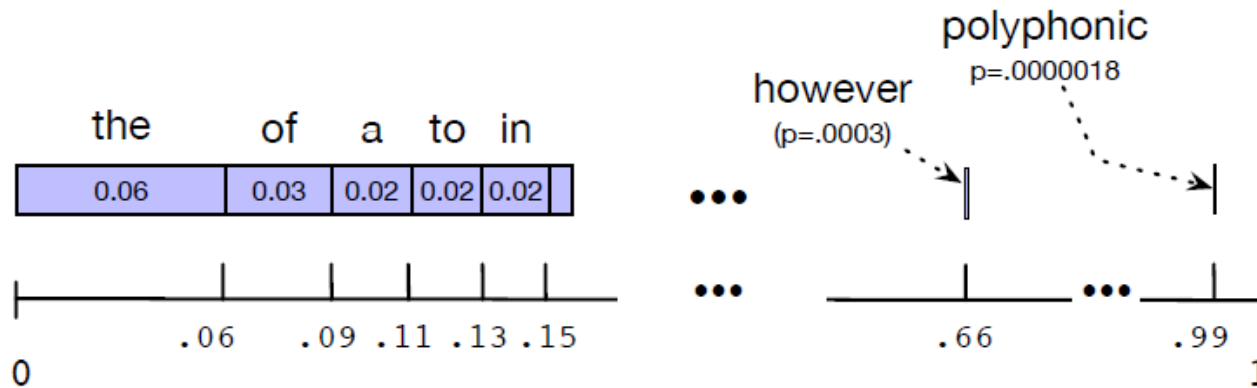


**Note:** $\quad 1 = \sum_i P(w_i)$

Randomly select value between 0 and 1: $\qquad$ e.g. 0.10

GENERATED SEQUENCE $\qquad$ => $\qquad$ a
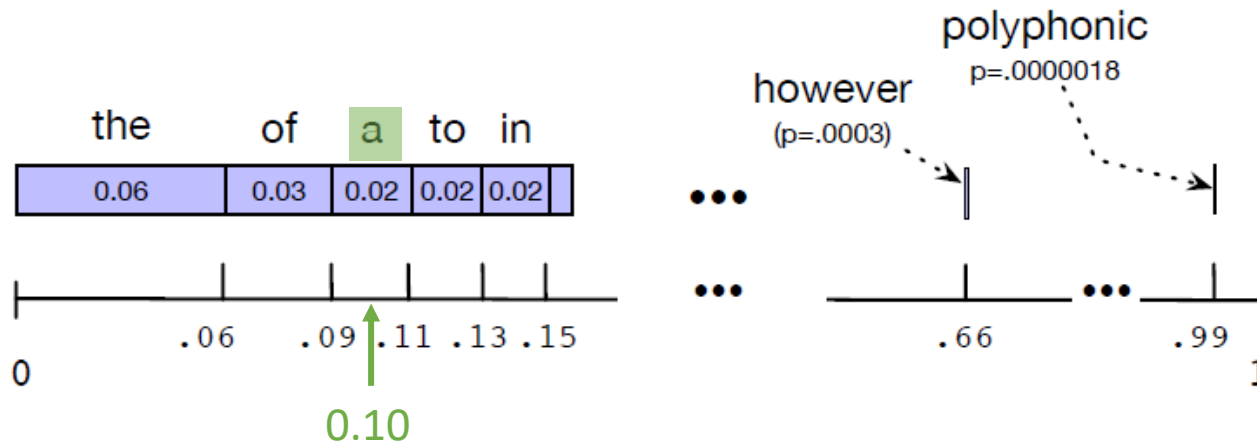
# Digression

## Sample sentence from a *unigram* LM

Assume we have access to words in $\quad V \cup \{< EOS >\}$

Assume sorted word probabilites: $\quad P(w_1) > P(w_2) > \ldots > P(w_{|V|+1})$



**Note:** $\quad 1 = \sum_i P(w_i)$

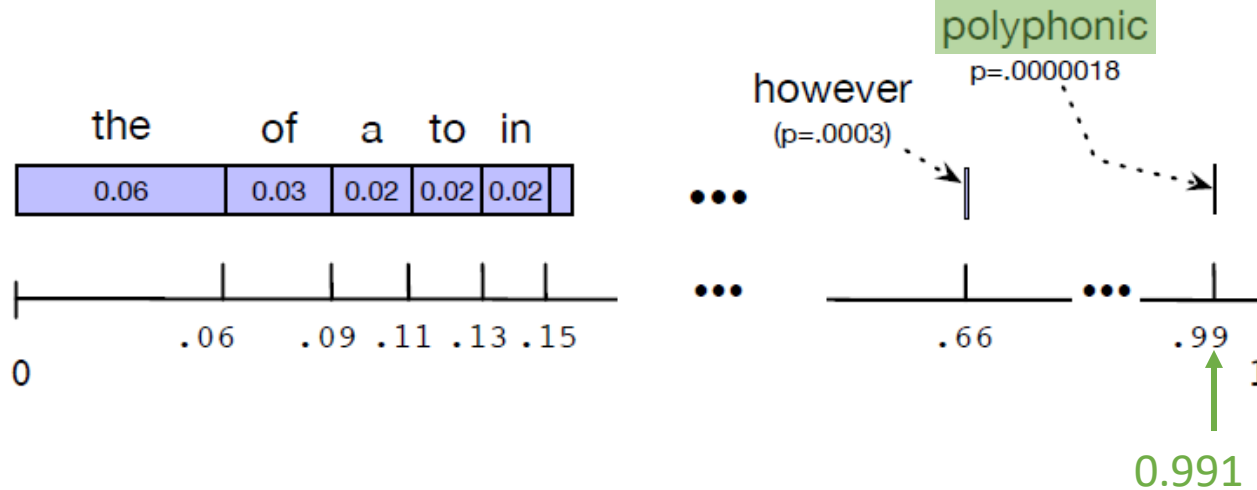Randomly select value between 0 and 1: $\qquad$ e.g. 0.991

GENERATED SEQUENCE $\qquad$ => $\qquad$ a polyphonic

Technische Hochschule
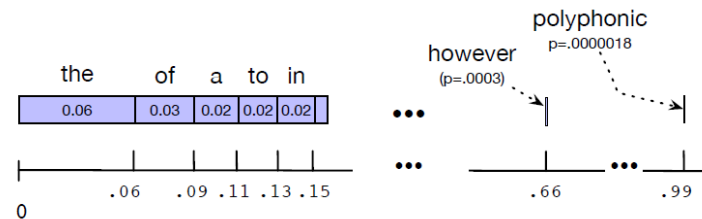Ingolstadt

# Digression
## Sample sentence from a *unigram* LM

Assume we have access to words in $\qquad V \cup \{< EOS >\}$

Assume sorted word probabilites: $\qquad P(w_1) > P(w_2) > \ldots > P\left(w_{|V|+1}\right)$



**Note:** $\quad 1 = \sum_i P(w_i)$

**Algorithm**

Repeat below steps 1-3, until you generate <EOS> token

1.  Randomly select a value $x \in [0,1]$
2.  Find that point $x$ on the line
3.  Print word whose interval includes chosen value $x$

Further reading: Section 3.3, https://web.stanford.edu/~jurafsky/slp3/3.pdf

# Further questions

[1] https://web.stanford.edu/~jurafsky/slp3/4.pdf

Technische **Hochschule**
**Ingolstadt**

# Further questions

- How to compute $P(c)$? (see Section 4.2 in [1])

$$\hat{c} = \text{argmax}_{c \in C} \log P(c) + \sum_{i \in X} \log P(w_i | c)$$

[1] https://web.stanford.edu/~jurafsky/slp3/4.pdf

# Further questions

- How to compute $P(c)$? (see Section 4.2 in [1])
- What about words in test data which were not present in training documents for a class? (see Section 4.2 in [1])

$$\hat{c} = \text{argmax}_{c \in C} \log P(c) + \sum_{i \in X} \log P(w_i|c)$$

$$\hat{P}(w_i|c) = \frac{count(w_i, c)}{\sum_{w \in V} count(w, c)}$$

[1] https://web.stanford.edu/~jurafsky/slp3/4.pdf

# Further questions

- How to compute $P(c)$? (see Section 4.2 in [1])

- What about words in test data which were not present in training documents for a class? (see Section 4.2 in [1])

- How to evaluate the classification algorithm?
  (next topic or Section 4.7 in [1])

[1] https://web.stanford.edu/~jurafsky/slp3/4.pdf