# COMP90042
# Web search and text analysis

## Workshop Week 6

xudong.han@unimelb.edu.au
https://github.com/HanXudong/COMP90042_Workshops

# Review

1. N-gram model

2. Backoff and interpolation

1. &lt;s&gt; &lt;s&gt; how much wood would a wood chuck chuck if a wood chuck would chuck wood &lt;/s&gt;

2. &lt;s&gt; &lt;s&gt; a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood <span style="color:red">&lt;/s&gt;</span>

&lt;s&gt; &lt;s&gt; a wood chuck would chuck the wood he could chuck if a wood chuck would chuck wood <span style="color:red">&lt;/s&gt;&lt;/s&gt;</span>

$$P(w_1, w_2, \ldots, w_m) = \prod_{i=1}^{m} P(w_i \mid w_{i-2} \ w_{i-1})$$

$$P_{add1}(w_i \mid w_{i-2} \ w_{i-1}) = \frac{C(w_{i-2} \ w_{i-1} \ w_i) + 1}{C(w_{i-2} \ w_{i-1}) + V}$$

# Q3: What does back–off mean, in the context of smoothing a language model? What does interpolation refer to?

- The idea in a Backoff model is to build an Ngram model based on an (N-1) model

- https://en.wikipedia.org/wiki/Katz%27s_back-off_model

- Interpolation: instead of just backing off to the non-zero Ngram, it is possible to take into account all Ngrams.

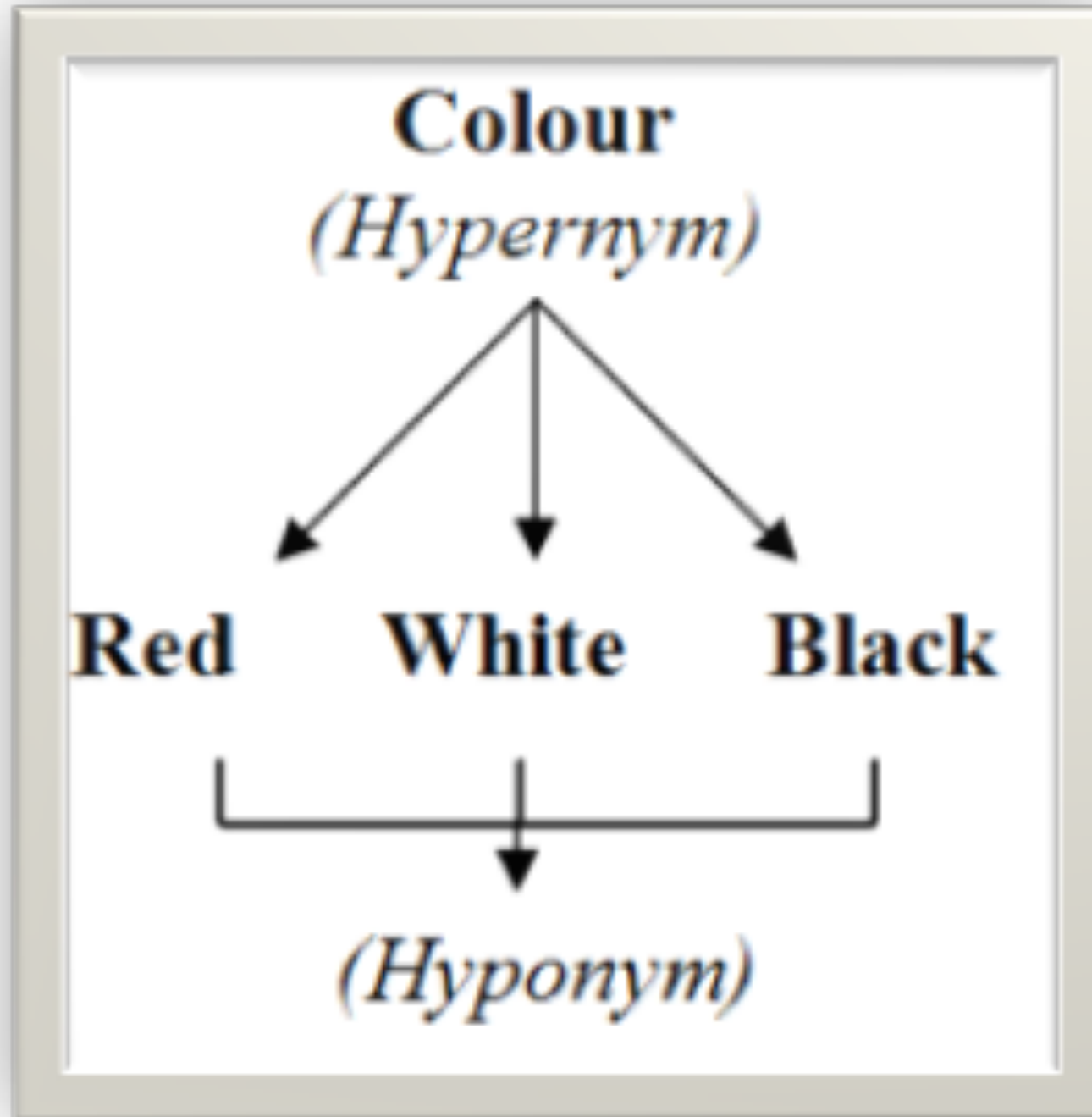- Estimate lambdas from held-out dataset.

# This workshop

- Words and senses

- Wordnet and lexical semantics

- Distributional semantics and word embedding

# Senses

# Senses

# Senses

- **Meronym:** Part of a whole



- **Holonym:** The whole to which parts belong

**Q2**

# WordNet

https://github.com/evanmiltenburg/images/tree/master/wordnet%20graphs

```
                                      entity
                                      abstraction...
                                      communication
                                      message...                                    entity
 entity               entity         statement           entity                     abstraction...
 abstraction...       abstraction... pleading            abstraction...             measure
 communication        psychological..charge...           group...                   system of meas...
 message...           cognition...   accusation...       collection...              information meas...
```
information

```
 entity               entity
 physical...          abstraction...
 process...           psychological...    entity
 processing           cognition...        abstraction...
 data process...      process...          psychological...
 operation            basic cog...        event
 computer op...       memory...           act...
```
retrieval

*information* is more similar to the word *retrieval* or the word *science*

$$WuP\_sim(c_1, c_2) = \frac{2 \times depth(LCS(c_1, c_2))}{depth(c_1) + depth(c_2)}$$

|  |  | information | | | | |
|---|---|---|---|---|---|---|
|  |  | 1 | 2 | 3 | 4 | 5 |
| retrieval | 1 | 0.154 | 0.154 | 0.118 | 0.154 | 0.143 |
|  | 2 | 0.308 | 0.615 | 0.235 | 0.308 | 0.286 |
|  | 3 | 0.364 | 0.545 | 0.267 | 0.364 | 0.333 |

# Q4a PMI

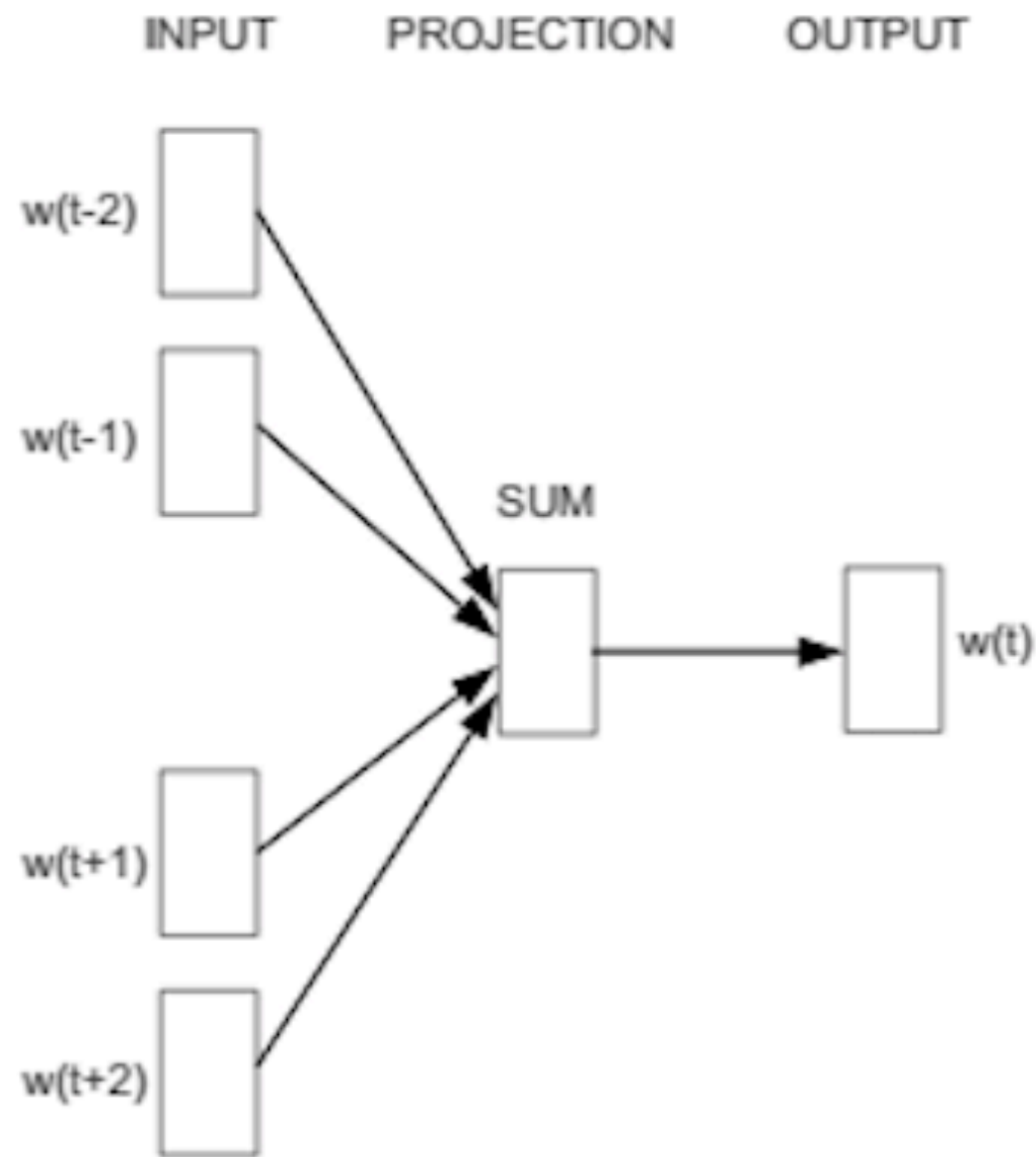| | cup | not (cup) | Total |
|---|---:|---:|---:|
| **world** | 55 | 225 | **280** |
| **not (world)** | 315 | 1405 | **1720** |
| **Total** | **370** | **1630** | **2000** |

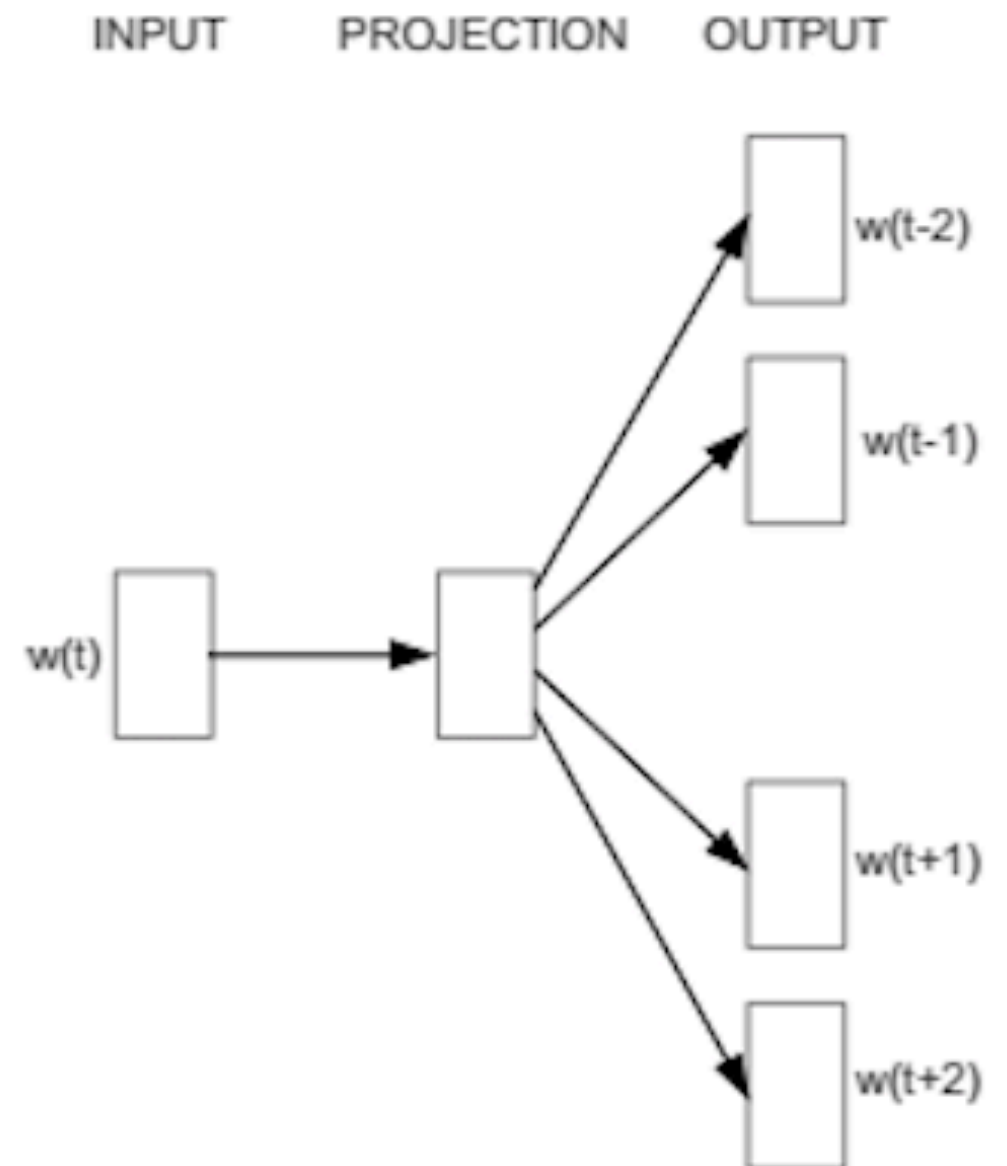$$PMI(x, y) = log_2 \frac{p(x, y)}{p(x)p(y)} = log_2 P(x, y) - log_2 p(x) - log_2 P(y)$$

$$PMI(x, y) = log_2 \frac{p(x, y)}{p(x)p(y)} = log_2 \frac{p(y \mid x)}{p(y)} = log_2 \frac{p(x \mid y)}{x}$$

# Q6 word to vector
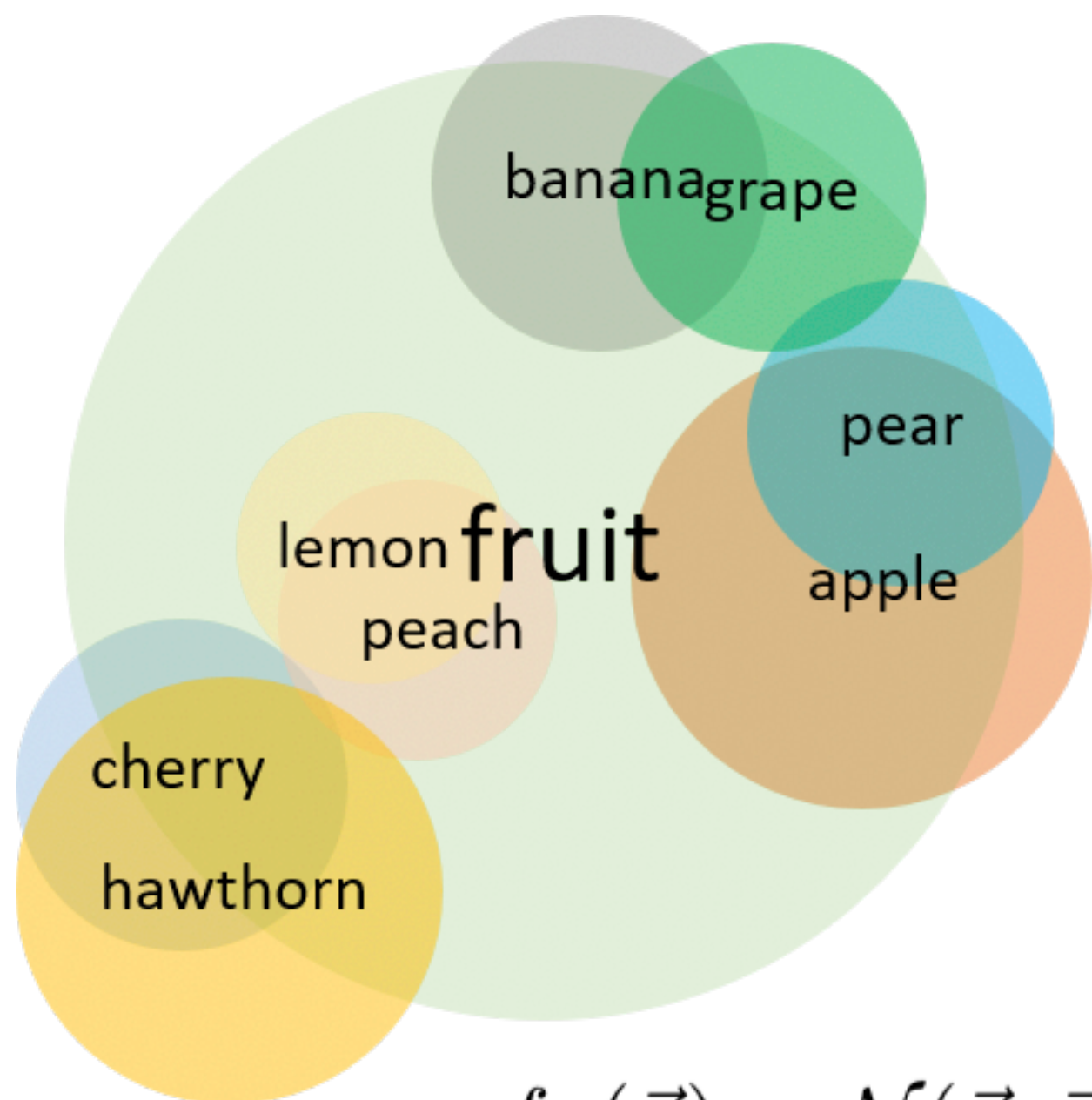
CBOW

Skip-gram

Gaussian Word Embedding

$$f_w(\vec{x}) = \mathcal{N}(\vec{x}; \vec{\mu}_w, \Sigma_w)$$

$$= \frac{1}{\sqrt{(2\pi)^D |\Sigma_w|}} e^{-\frac{1}{2}(\vec{x} - \vec{\mu_w})^\top \Sigma_w^{-1} (\vec{x} - \vec{\mu_w})}$$