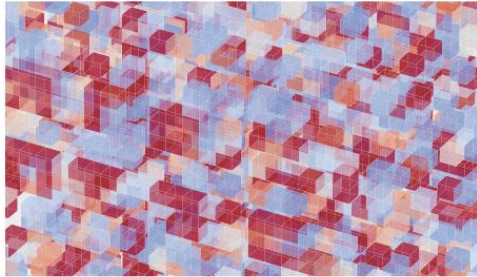# The Structure of Meaning in Language

# Why can computers understand natural language

# The Structure of Meaning in Language: Parallel Narratives in Linear Algebra and Category Theory

*Tai-Danae Bradley, Juan Luis Gastaldi, and John Terilla*

## Introduction

Categories for AI, an online program about category theory in machine learning, unfolded over several months beginning in the fall of 2022. As described on their website `https://cats.for.ai`, the "Cats for AI" organizing committee, which included several researchers from industry including two from DeepMind, felt that the machine learning community ought to be using more rigorous compositional tools and that category theory has "great potential to be a cohesive force" in science in general and in artificial intelligence in particular. While this article is by no means a comprehensive report on that event, the popularity of "Cats for AI" — the five introductory lectures have been viewed thousands of times — signals the growing prevalence of category theoretic tools in AI.

One way that category theory is gaining traction in machine learning is by providing a formal way to discuss how learning systems can be put together. This article has a different and somewhat narrow focus. It's about how a fundamental piece of AI technology used in language modeling can be understood, with the aid of categorical thinking, as a process that extracts structural features of language from purely syntactical input. The idea that structure arises from form may not be a surprise for many readers — category theoretic ideas have been a major influence in pure mathematics for three generations — but there are consequences for linguistics that are relevant for some of the ongoing debates about artificial intelligence. We include a section that argues that the mathematics in these pages rebut some widely accepted ideas in contemporary linguistic thought and support a return to a structuralist approach to language.

The article begins with a fairly pedantic review of linear algebra which sets up a striking parallel with the relevant category theory. The linear algebra is then used to review how to understand word embeddings, which are at the root of current large language models (LLMs). When the linear algebra is replaced, *Mad Libs* style, with the relevant category theory, the output becomes not word embeddings but a lattice of formal concepts. The category theory that gives rise to the concept lattice is a particularly simplified piece of enriched category theory and suggests that by simplifying a little less, even more of the structure of language could be revealed.

*Tai-Danae Bradley is a research mathematician at SandboxAQ and a visiting faculty member at The Master's University. Her email address is* tai.danae@math3ma.com.

*Juan Luis Gastaldi is a researcher at ETH Zürich. His email address is* juan.luis.gastaldi@gess.ethz.ch.

*John Terilla is a professor of mathematics at CUNY Queens College and on the Doctoral Faculty at the CUNY Graduate Center. His email address is* jterilla@gc.cuny.edu.

---

# Why can computers understand natural language?

### The structuralist image of language behind word embeddings
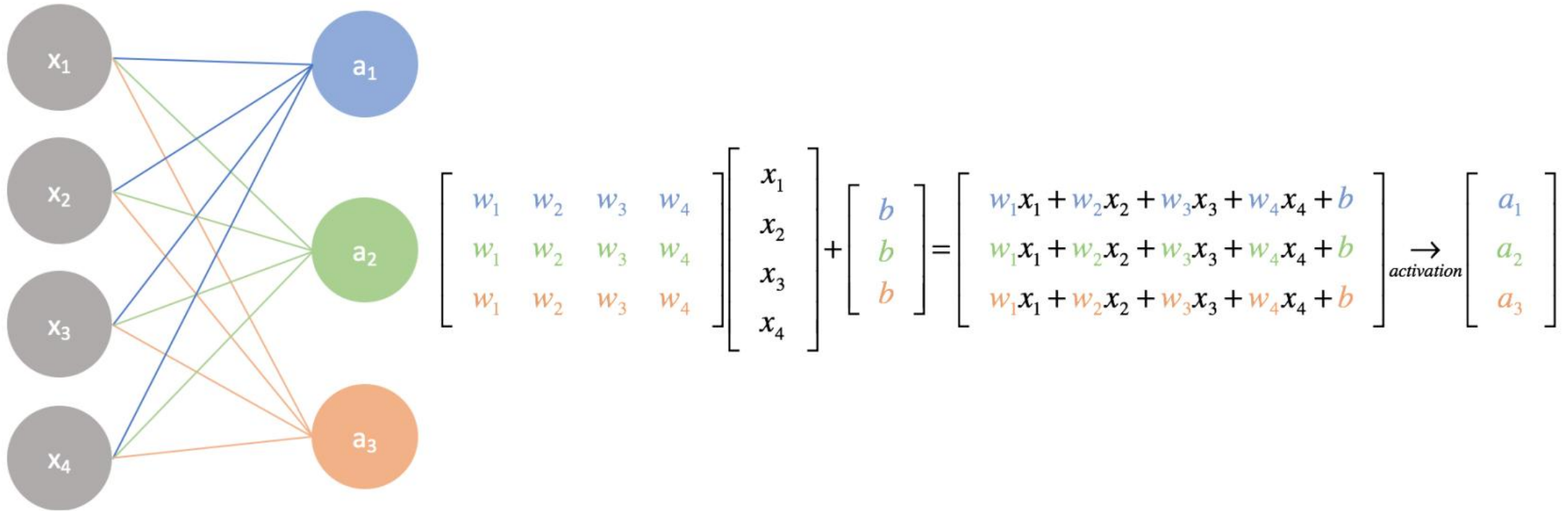
Juan Luis Gastaldi

#### Abstract

The present paper intends to draw the conception of language implied in the technique of word embeddings that supported the recent development of deep neural network models in computational linguistics. After a preliminary presentation of the basic functioning of elementary artificial neural networks, we introduce the motivations and capabilities of word embeddings through one of its pioneering models, word2vec. To assess the remarkable results of the latter, we inspect the nature of its underlying mechanisms, which have been characterized as the implicit factorization of a word-context matrix. We then discuss the ordinary association of the "distributional hypothesis" with a "use theory of meaning", often justifying the theoretical basis of word embeddings, and contrast them to the theory of meaning stemming from those mechanisms through the lens of matrix models (such as VSMs and DSMs). Finally, we trace back the principles of their possible consistency through Harris's original distributionalism up to the structuralist conception of language of Saussure and Hjelmslev. Other than giving access to the technical literature and state of the art in the field of Natural Language Processing to non-specialist readers, the paper seeks to reveal the conceptual and philosophical stakes involved in the recent application of new neural network techniques to the computational treatment of language.
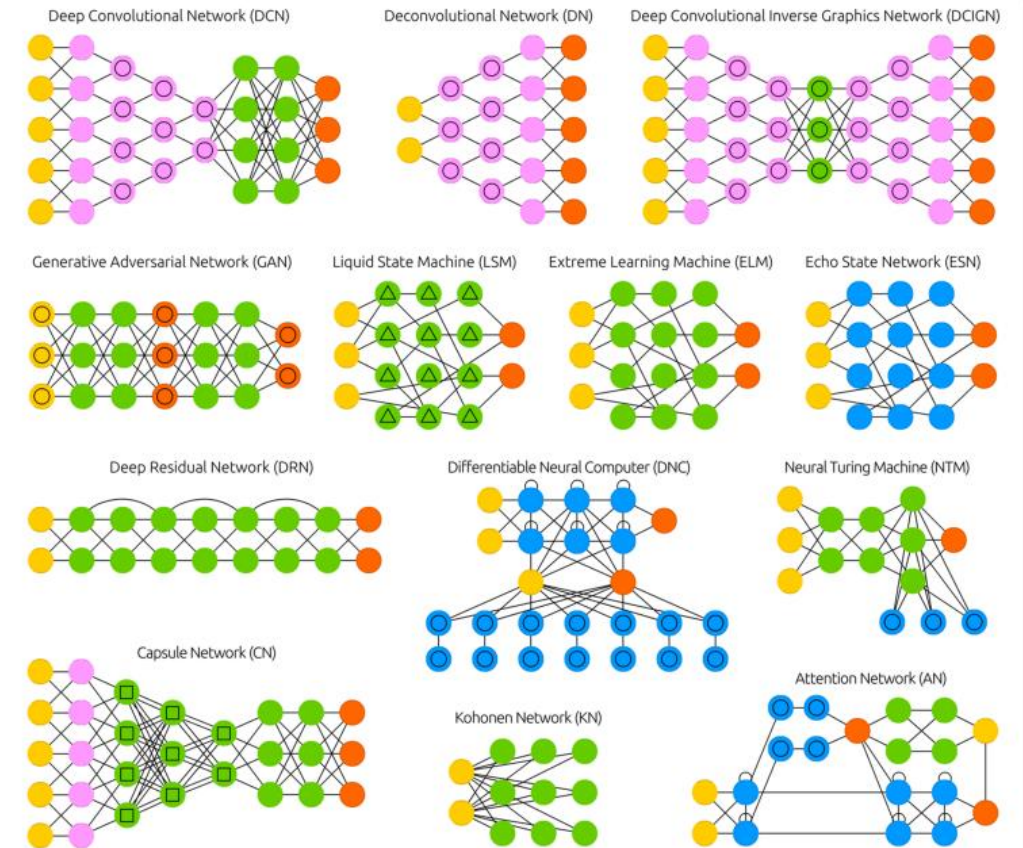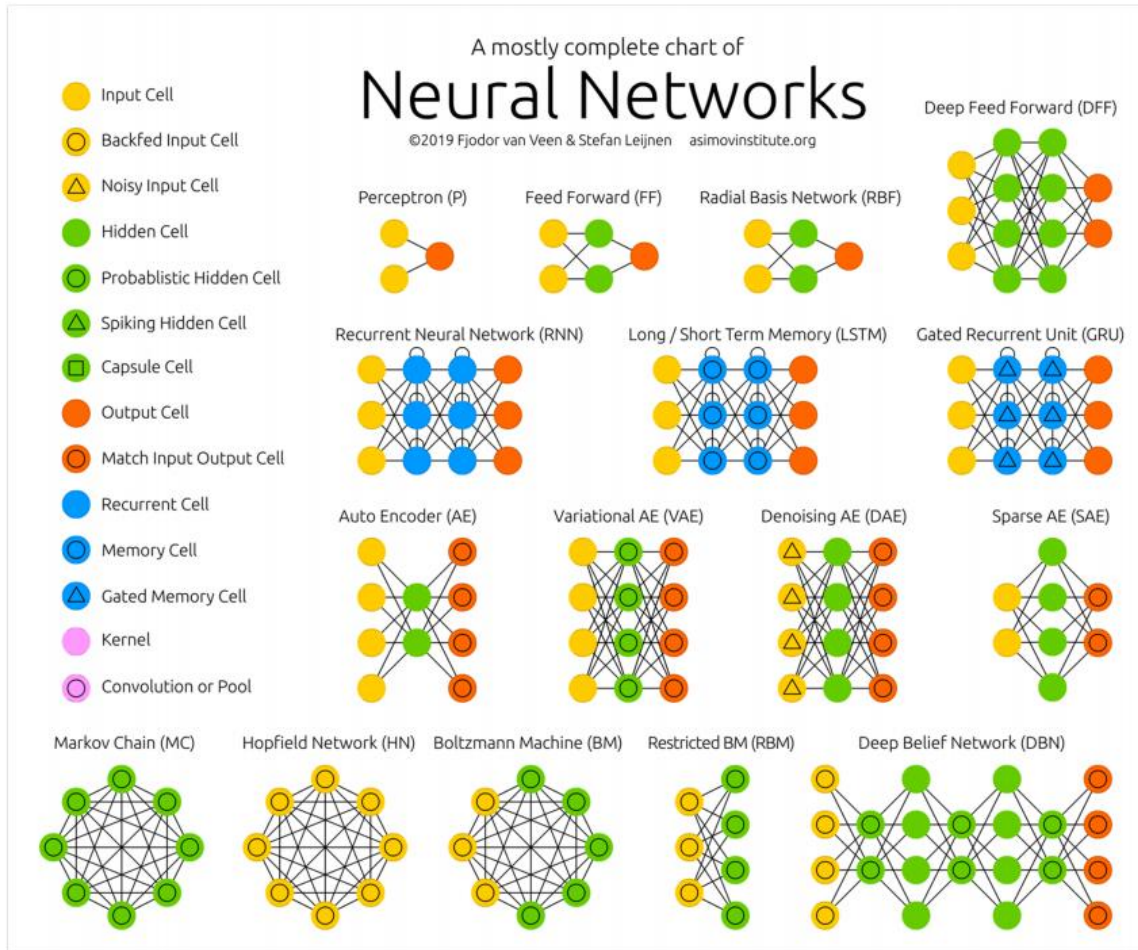
**Keywords** Word Embeddings · Natural Language Processing · word2vec · Neural Networks · Philosophy of Language · Matrix Models · Distributional Hypothesis · Structuralism

# Neural Networks



$$\begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b \\ b \\ b \end{bmatrix} = \begin{bmatrix} w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b \\ w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b \\ w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b \end{bmatrix} \xrightarrow[activation]{} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

# The Family of DNNs

# Deep Neural Nets (DNNs)

$$x \quad = \quad z^{(1)} = a^{(1)}$$

$$z^{(2)} \quad = \quad W^{(1)}x + b^{(1)}$$

$$a^{(2)} \quad = \quad f\left(z^{(2)}\right)$$

$$z^{(3)} \quad = \quad W^{(2)}a^{(2)} + b^{(2)}$$

$$a^{(3)} \quad = \quad f\left(z^{(3)}\right)$$

$$s \quad = \quad U^T a^{(3)}$$

# DNNs and Natural Language I

| Index | Word |
|-------|------|
| … | … |
| 535 | nearly |
| 536 | shares |
| 537 | member |
| 538 | campaign |
| 539 | media |
| 540 | needs |
| 541 | why |
| 542 | house |
| 543 | issues |
| 544 | costs |
| 545 | fire |
| … | … |

$$v_{house} = (0,0,0,0,0,0,0,0,\ldots,0,\overbrace{1}^{542^{nd}\ position},0,\ldots,\underbrace{0,0,0,0,0,0,0,0)}_{3\ million\ dimensions}$$

One-hot coding

# Word Embeddings: word2vec

# Dense Vector Representations

30, 000-dimensional real vector space

$$v_{house} = (0, 0, 0, 0, 0, 0, 0, 0, \ldots, \overbrace{1}^{542^{nd} \text{ position}}, 0, \ldots, 0, 0, 0, 0, 0, 0, 0, 0)$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{\text{million dimensions}}$$

$$v_{house} = (\underbrace{0.157227, -0.0708008, 0.0539551, \ldots, -0.041748, 0.00982666, -0.00494385, -0.032959}_{300 \text{ dimensions}})$$

$$f : \boxed{\mathbb{R}^n \xrightarrow{f_1} \mathbb{R}^{n_1}} \xrightarrow{f_2} \mathbb{R}^{n_2} \xrightarrow{f_3} \ldots$$

$$\ldots \xrightarrow{f_K} \mathbb{R}^{n_K} \xrightarrow{g} \mathbb{R}^m \qquad (1)$$

$$f_i(\mathbf{x}) = a(\mathbf{M}_i \mathbf{x} + b_i), \qquad (2)$$

# meaning of words

- Not only did the performance of models across different tasks increase substantially,

- but also unexpected linguistic significance was found in the vector space operations of the embedded word vectors. In particular, the inner product between two vectors shows a high correlation with semantic similarity
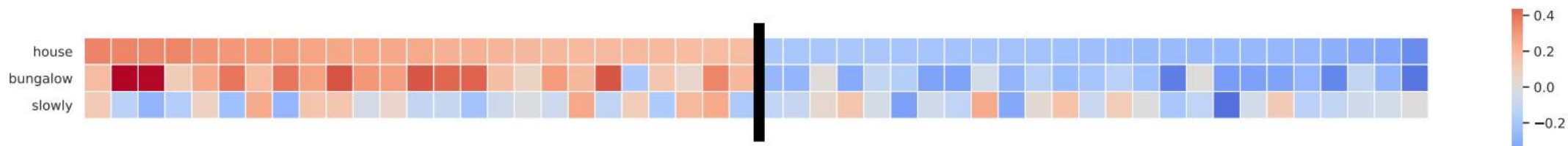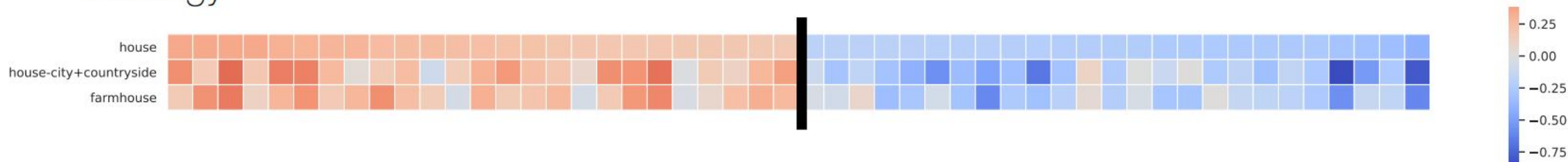
◇ Example: `house`


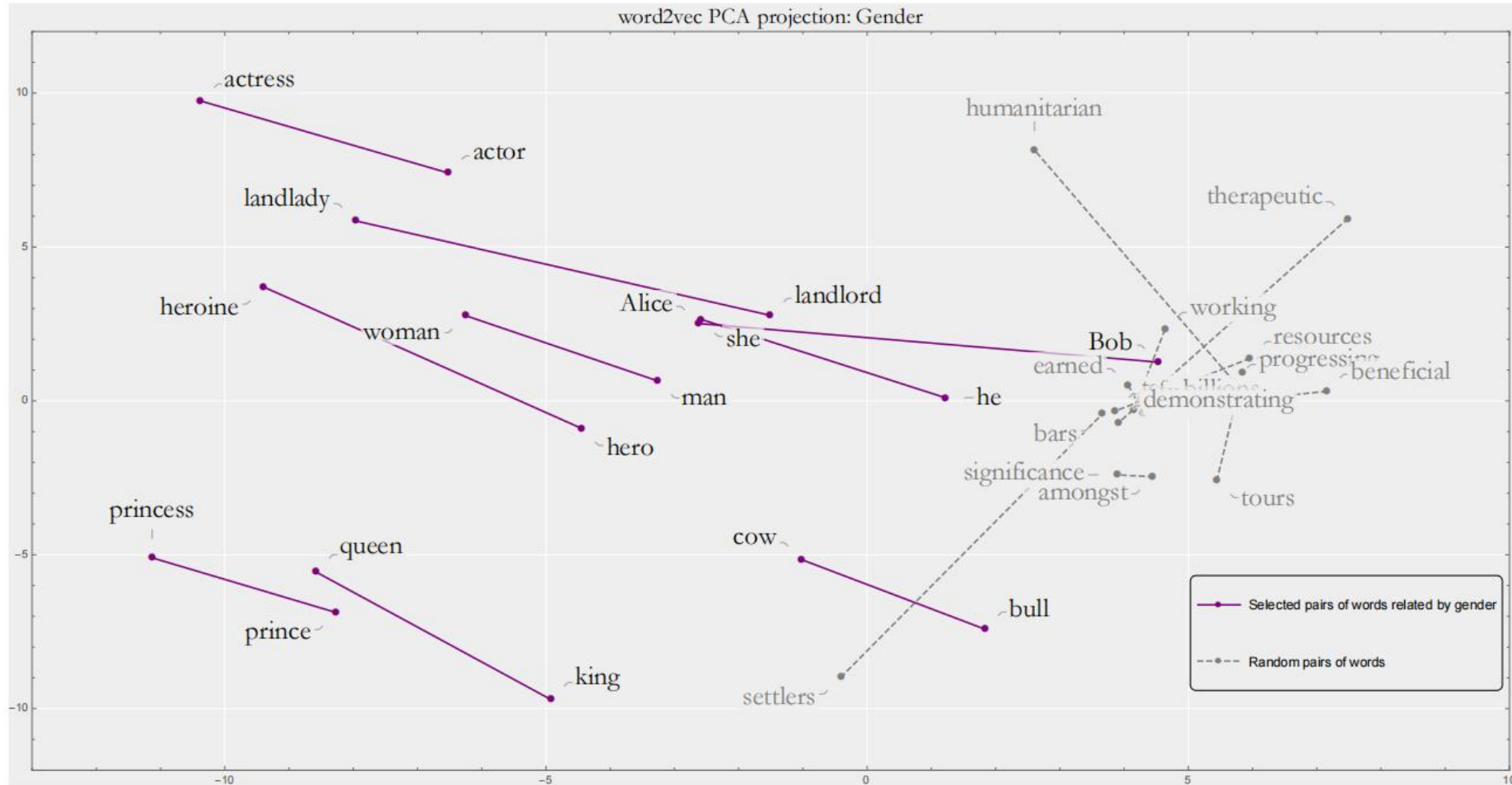
◇ Syntactic and semantic properties

   &mdash; Similarity



   &mdash; Analogy

# Word Embeddings: Similarity

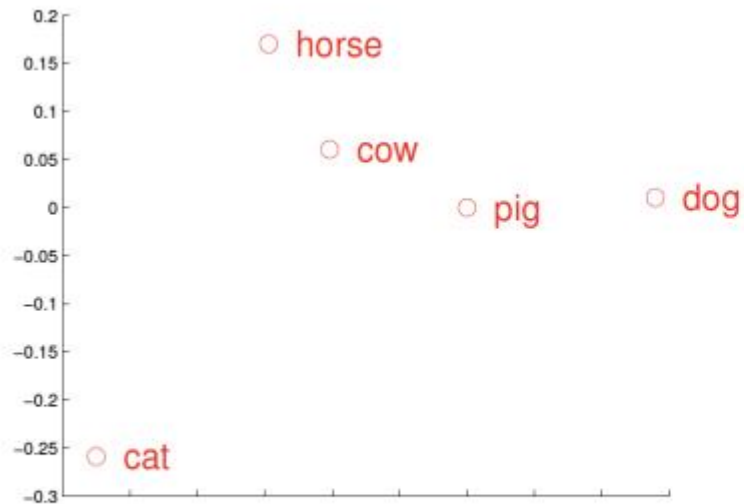| house | cosine distance |
|---|---|
| houses | 0.292761 |
| bungalow | 0.312144 |
| apartment | 0.3371 |
| bedroom | 0.350306 |
| townhouse | 0.361592 |
| residence | 0.380158 |
| mansion | 0.394181 |
| farmhouse | 0.414243 |
| duplex | 0.424206 |
| homes | 0.43802 |



(https://projector.tensorflow.org)

# Word Embeddings: Analogy

$$v_{king} - v_{queen} \approx v_{hero} - v_{heroine}$$
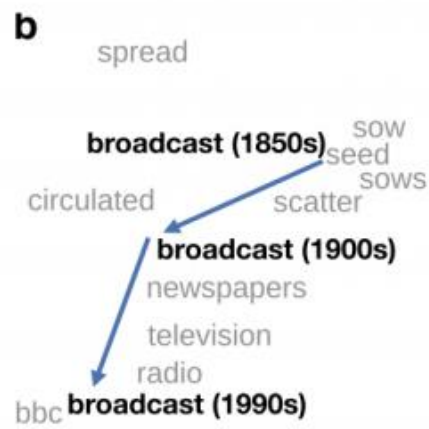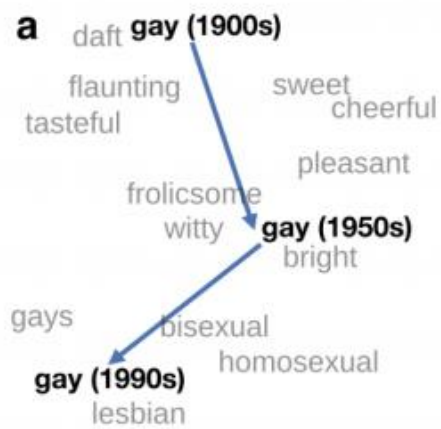


word2vec PCA projection: Gender

# Word Embeddings: Analogy

$$v_{good} - v_{better} \approx v_{soft} - v_{softer}$$



word2vec PCA projection: Comparatives

(Mikolov, Sutskever, et al., 2013)



(Hamilton et al., 2016)

How to understand one-hot coding ?

There is a natural embedding

$$\text{Set} \longrightarrow \text{functions on Set}$$

$$X \longrightarrow \text{Fun}(X)$$

$$x \longmapsto \delta(x,-)$$

If $k$ is a field, then $k^X$ is a Hilbert space

linear space, inner product. extra structure.

If $k$ is $\{0,1\}$, then $\{0,1\}^X$ can regard as
the set of all subset of $X$. $\{0,1\}^X$ is a Lattice
Boolean algebra.

## Category analogue

$C$ is a Category. we have Yoneda embedding 米田

$$y : C \longrightarrow Set^{C^{op}} = Fun(C^{op}, Set)$$

$$X \longmapsto hom(-, X)$$

$Set^{C^{op}}$ has all colimit.

And co-Yoneda embedding

$$z : C \longrightarrow (Set^C)^{op}$$

$$X \longmapsto hom(X, -)$$

$(Set^C)^{op}$ has all limit.

# Singular Value Decomposition (SVD)

$$M = U\Sigma V^*$$

Where:

$$
\begin{aligned}
M \quad &= \quad m \times n \text{ (real or complex) matrix} \\
U \quad &= \quad m \times m \text{ unitary matrix} \\
\Sigma \quad &= \quad m \times n \text{ non-negative real rectangular diagonal matrix} \\
V^* \quad &= \quad \text{conjugate transpose of } V, \text{ a } n \times n \text{ unitary matrix}
\end{aligned}
$$

In particular:

◇ The columns of $U$ (left singular vectors) are eigenvectors of $M \times M^*$

◇ The rows of $V^*$ (right singular values) are eigenvectors of $M^* \times M$

◇ The non-zero elements of $\Sigma$ (non-zero singular values) are the square roots of the non-zero eigenvalues of $M \times M^*$ or $M^* \times M$

How to get SVD :

$$M = U \Sigma V^*$$

$$M^* = V \Sigma U^*$$

$$M^*M = V \Sigma U^* \cdot U \Sigma V^* = V \Sigma^2 V^*$$

$$MM^* = U \Sigma V^* \cdot V \Sigma U^* = U \Sigma^2 U^*$$
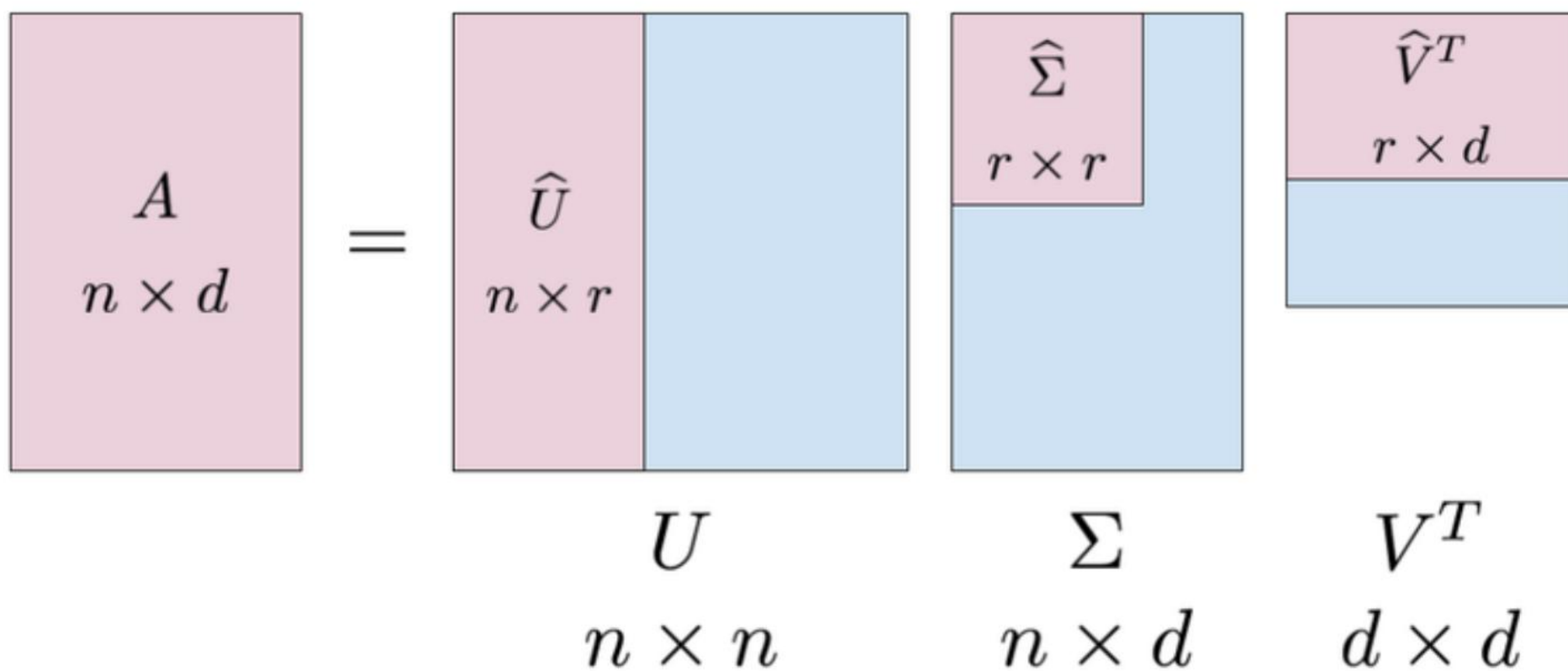
$$(M^*M)V = \Lambda V \qquad \sigma_i = \sqrt{\lambda_i}$$

$$(MM^*)U = \Lambda U$$

$$M^* u_j = \sigma_j v_j \qquad M v_j = \sigma_j u_j$$

$$M = U\Sigma V^*$$



$$A$$
$$n \times d$$

$$=$$

$$\widehat{U}$$
$$n \times r$$

$$\widehat{\Sigma}$$
$$r \times r$$

$$\widehat{V}^T$$
$$r \times d$$

$$U$$
$$n \times n$$

$$\Sigma$$
$$n \times d$$

$$V^T$$
$$d \times d$$

$$X = \{-, /, 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, =, a, b, c, d, e, f,$$
$$g, h, i, j, k, l, m, n, o, p, q, r, s, t, u, v, w, x, y, z, é\}.$$

$$Y = X \times X$$

the: $\quad h \longleftrightarrow t\_e$

hat: $\quad a \longleftrightarrow h\_t$

$\underset{\text{context}}{c} \in X \qquad\qquad \in X \times X = Y$

Linguistically relevant measure

$$m : \quad X \times Y \longrightarrow \mathbb{R}$$

$$(w, c) \longmapsto \log\left(\frac{p(w,c)}{p(w)p(c)}\right)$$

$$X \longrightarrow k^Y$$
$$x \longmapsto m(x, -)$$

$$Y \longrightarrow k^X$$
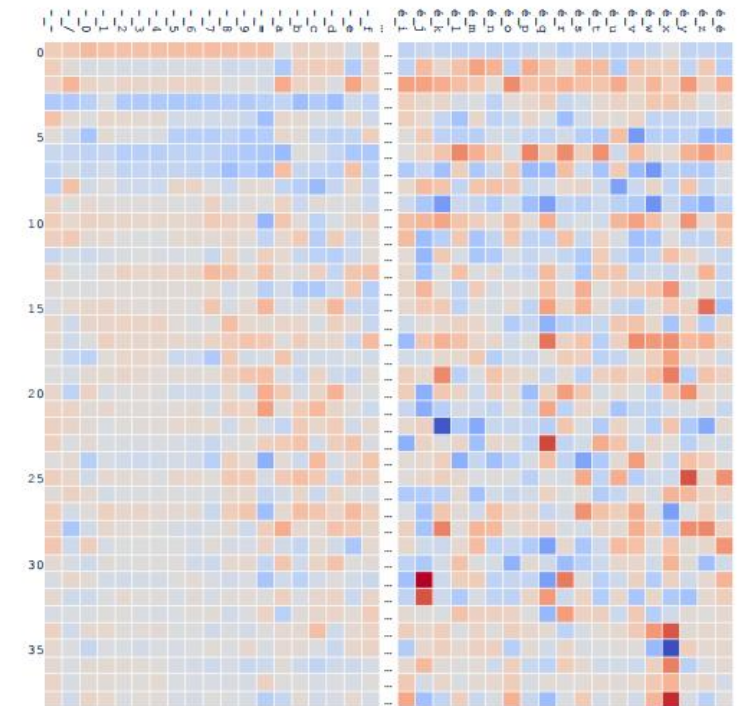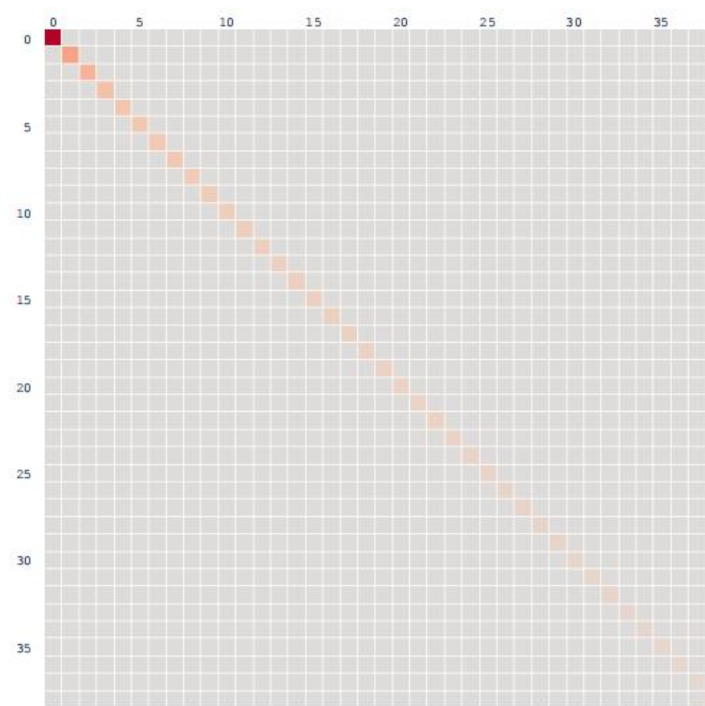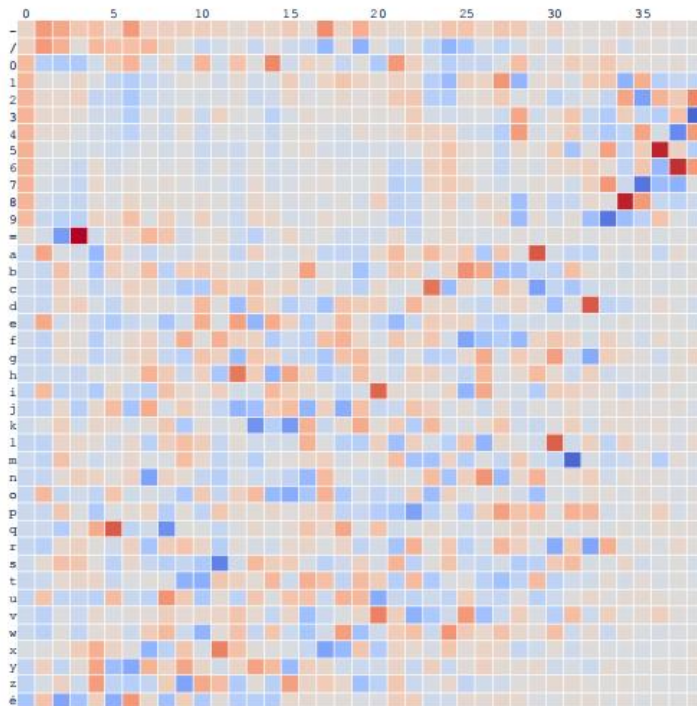$$y \longmapsto m(-, y)$$

---

$$k^Y$$

$$\uparrow \exists M^*$$

$$X \hookrightarrow k^X$$
$$x \longmapsto \delta(x, -)$$

$$\delta(-, y) \longleftarrow y$$

$$k^Y \longleftarrow Y$$

$$M \downarrow$$
$$k^X$$

$$MM^* : k^X \longrightarrow k^X$$
$$M^*M : k^Y \longrightarrow k^Y$$

# word2vec as Implicit Matrix Factorization
(Levy and Goldberg, 2014)

# Embeddings as Truncated SVD



$$\hat{V}^*$$

$$\hat{U} \times \hat{\Sigma}$$

$$M =$$
$$PMI(w,c) - \log k$$

$$M \approx$$
$$\hat{U} \times \hat{\Sigma} \times \hat{V}^*$$

$$PMI(w,c) =$$
$$\log \frac{p(w,c)}{p(w)p(c)}$$

# Example: Characters in Wikipedia



$$PMI(w,c) =$$

$$\log \frac{p(w,c)}{p(w)p(c)}$$
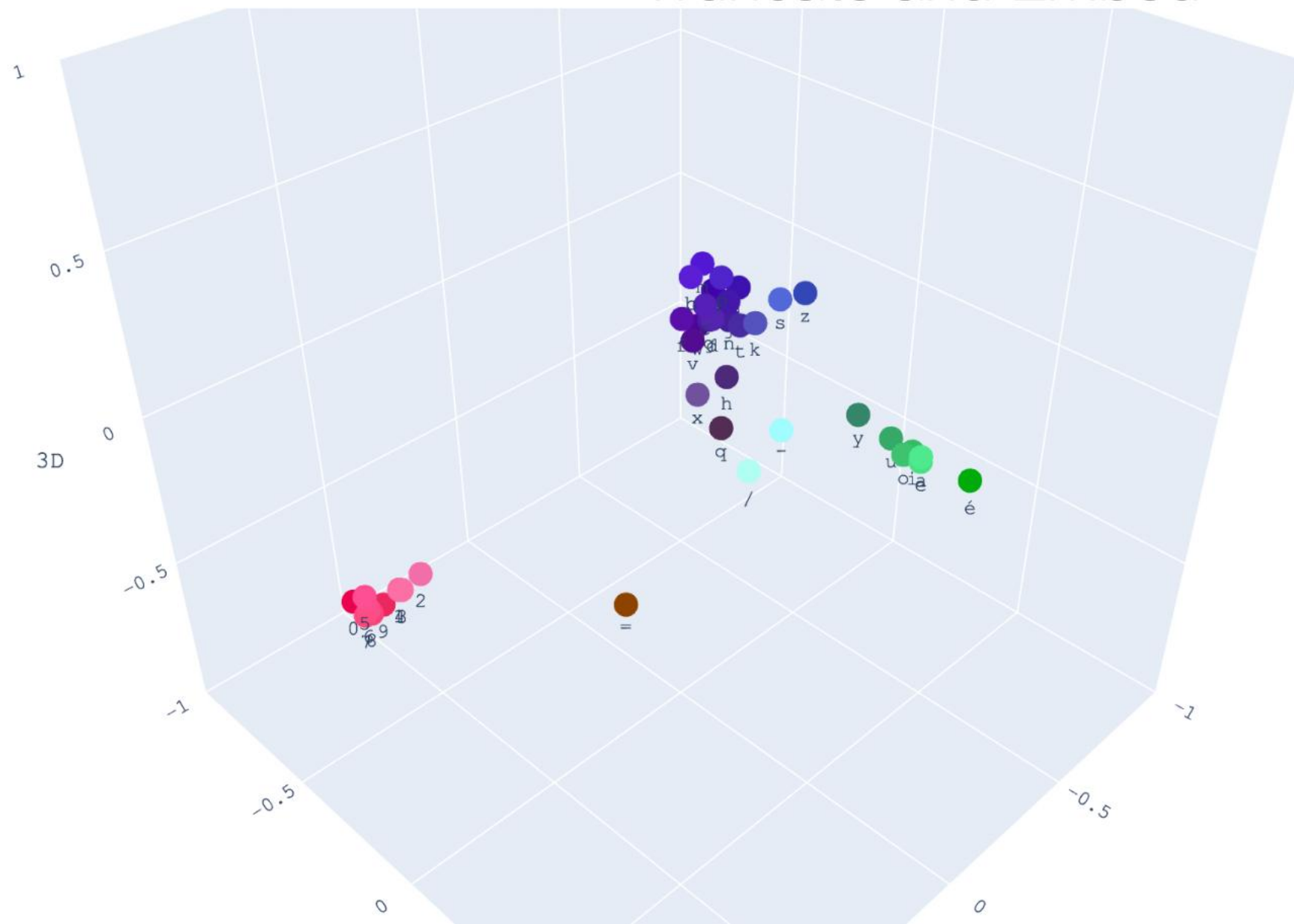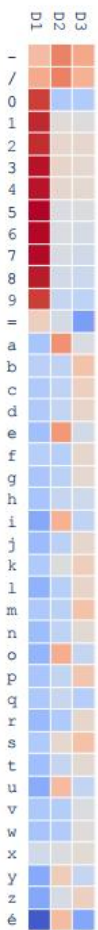
# SVD of Wikipedia Character PMI Matrix

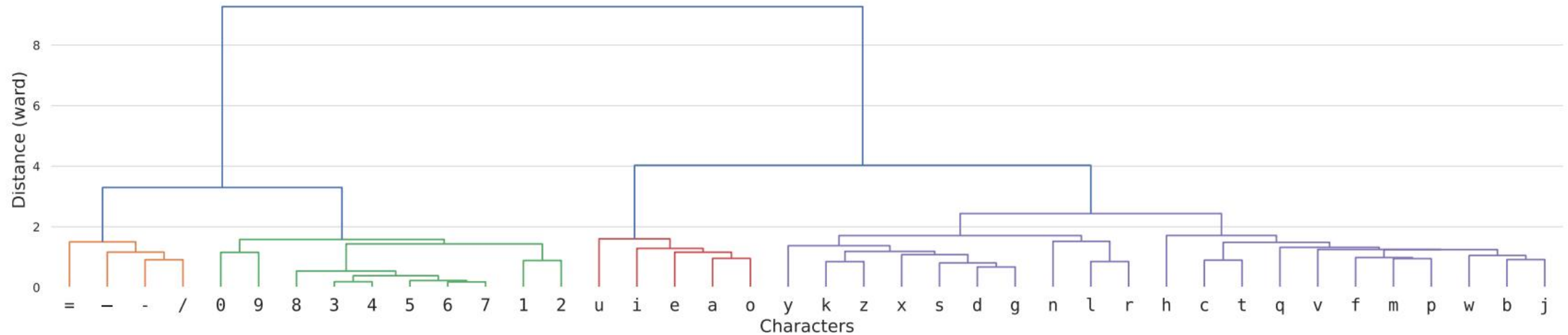# Every singular vectors represent some meaning

$$U \times \Sigma$$

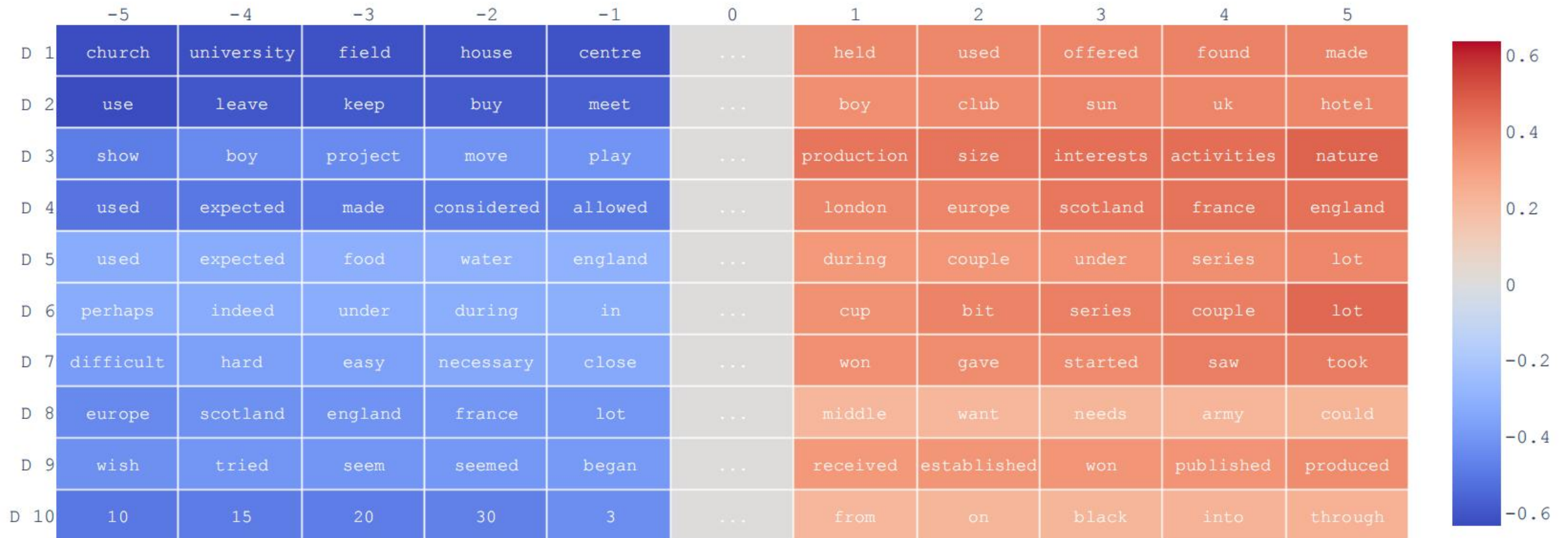Truncate and Embed

$\hat{U} \times \hat{\Sigma}$

# Clustering



$$O := \{=, -, \text{-}, /\}$$
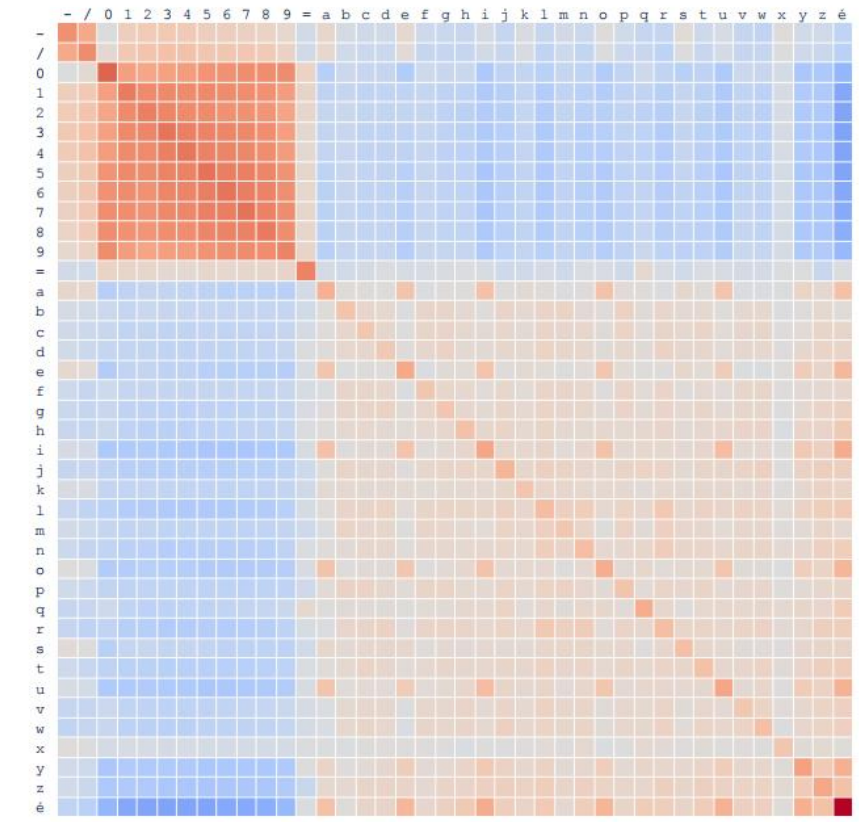
$$D := \{0, 9, 8, 3, 4, 5, 6, 7, 1, 2\}$$

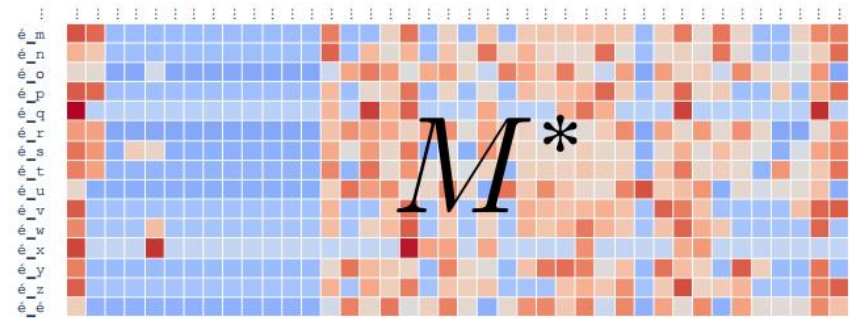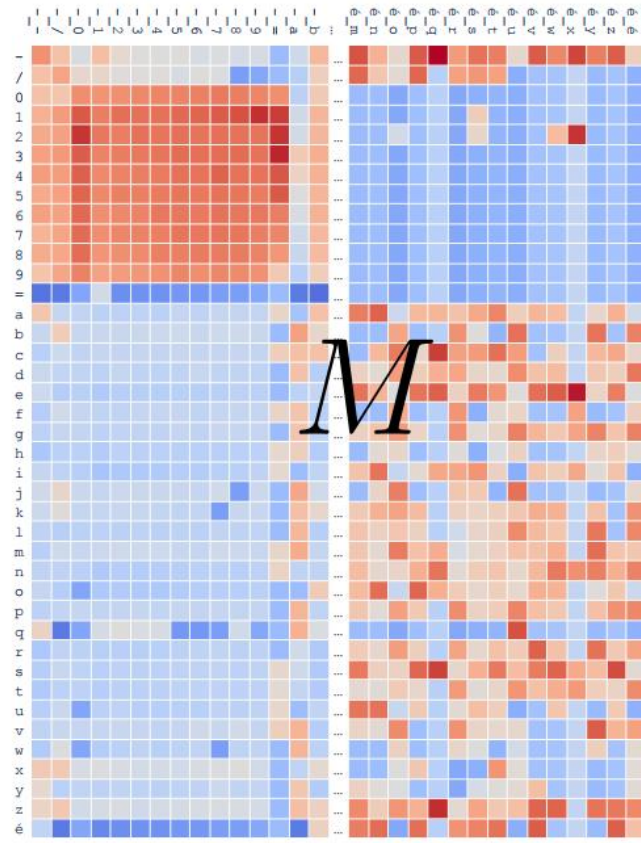$$V := \{u, i, e, a, o\}$$

$$C := \{y, k, z, x, s, d, g, n, l, r, h, c, t, q, v, f, m, p, w, b, j\}$$

# Words

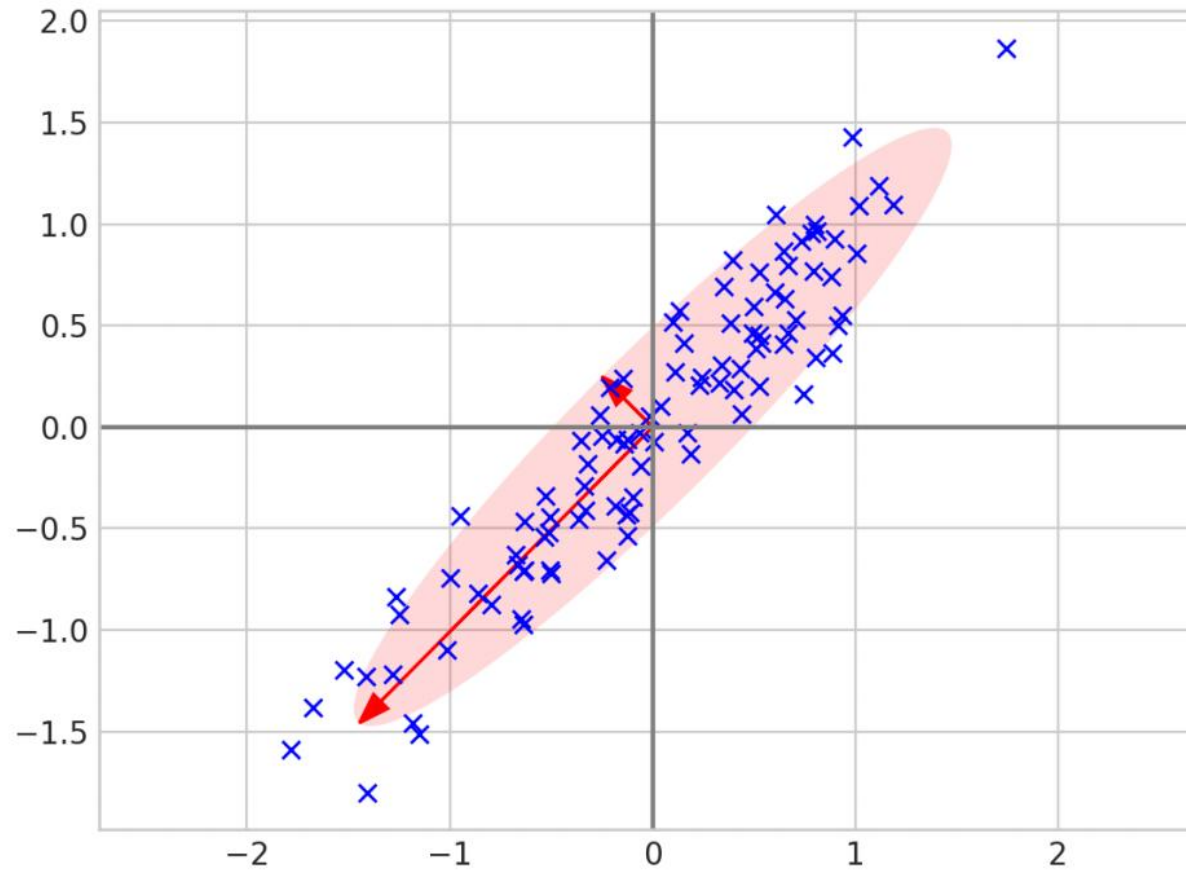| | -5 | -4 | -3 | -2 | -1 | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| D 1 | church | university | field | house | centre | ... | held | used | offered | found | made |
| D 2 | use | leave | keep | buy | meet | ... | boy | club | sun | uk | hotel |
| D 3 | show | boy | project | move | play | ... | production | size | interests | activities | nature |
| D 4 | used | expected | made | considered | allowed | ... | london | europe | scotland | france | england |
| D 5 | used | expected | food | water | england | ... | during | couple | under | series | lot |
| D 6 | perhaps | indeed | under | during | in | ... | cup | bit | series | couple | lot |
| D 7 | difficult | hard | easy | necessary | close | ... | won | gave | started | saw | took |
| D 8 | europe | scotland | england | france | lot | ... | middle | want | needs | army | could |
| D 9 | wish | tried | seem | seemed | began | ... | received | established | won | published | produced |
| D 10 | 10 | 15 | 20 | 30 | 3 | ... | from | on | black | into | through |

0.6
0.4
0.2
0
-0.2
-0.4
-0.6
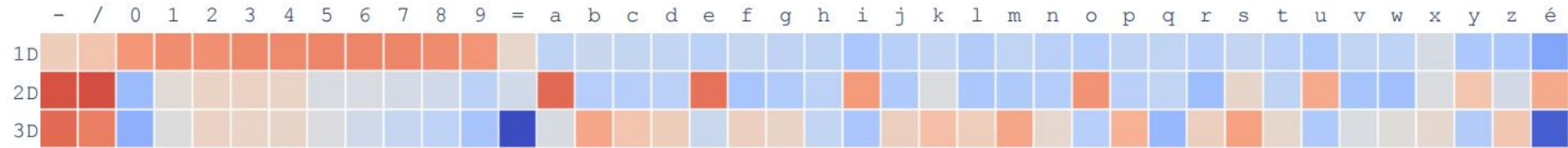
# $M \times M^*$ as A Covariance Matrix
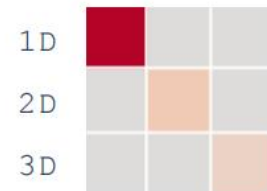
# Eigenvectors and Eigenvalues

# "Eigenstructure"

Eigenvectors of $M \times M^*$:
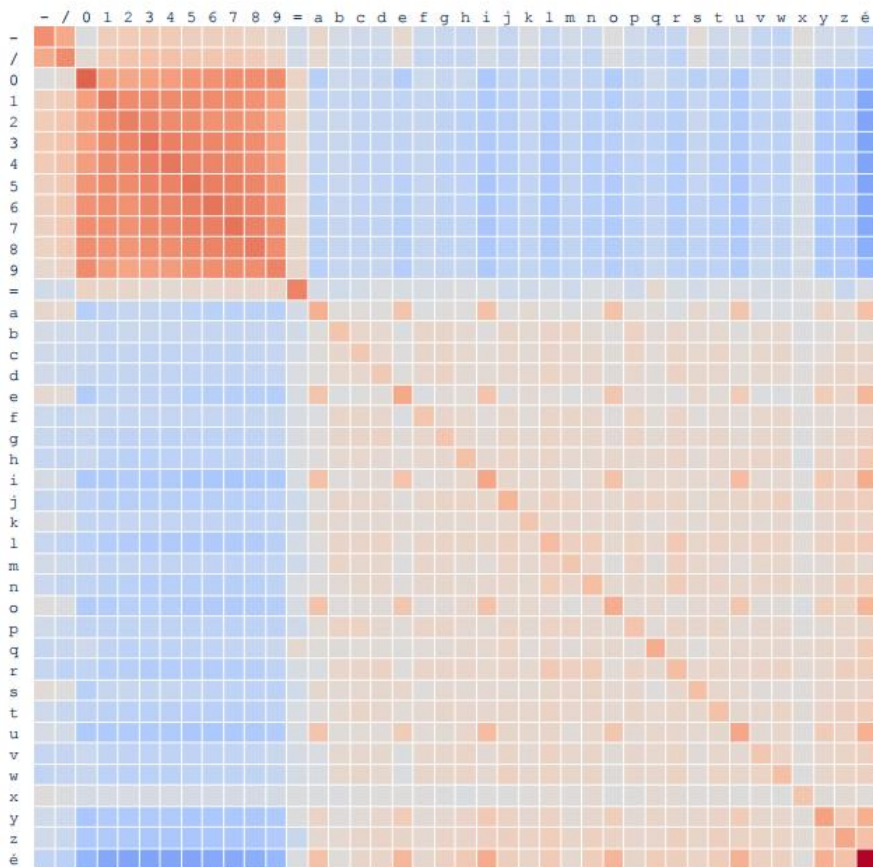


Eigenvalues of $M \times M^*$:



Eigenvectors of $M^* \times M$:
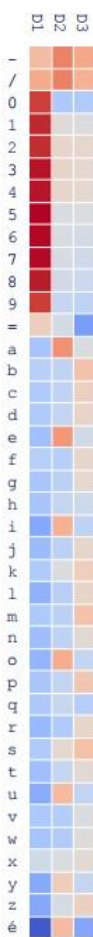
# Eigenvectors as Fixed Points

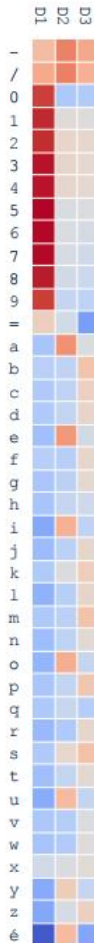$$(M \times M^*)v = \lambda v$$

Category analogue

$C$ is a category. we have Yoneda embedding 米田

$$y : C \longrightarrow Set^{C^{op}} = Fun(C^{op}, Set)$$

$$X \longmapsto hom(-, X)$$

$Set^{C^{op}}$ has all colimit.

And co-Yoneda embedding

$$z : C \longrightarrow (Set^{C})^{op}$$

$$X \longmapsto hom(X, -)$$

$(Set^{C})^{op}$ has all limit.

# Isbell Conjugation.

$$\mathrm{Fun}(C^{op}, \mathrm{Set}) \underset{F_*}{\overset{F^*}{\rightleftarrows}} \mathrm{Fun}(C, \mathrm{Set})^{op}$$

$$F^*: \mathrm{Set}^{C^{op}} \rightleftarrows (\mathrm{Set}^C)^{op} : F_*$$

$$F^*(f)(X) = \hom(f, y(x))$$

$$F_*(g)(X) = \hom(z(X), g)$$

profunctor : $m : C^{op} \times D \to Set$

$$C \to (Set^D)^{op}$$
$$c \mapsto m(c, -)$$

$$D \to Set^{C^{op}}$$
$$d \mapsto m(-, d)$$

$$C \to (Set^D)^{op}$$

$$(Set^D)^{op} \xleftarrow{\text{Coyoneda}} D$$

$$C \xrightarrow{\text{Yoneda}} Set^{C^{op}} \xrightarrow{F^*}$$

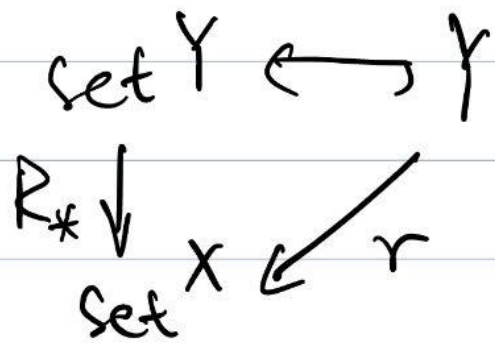$$F_* : (Set^D)^{op} \dashrightarrow Set^{C^{op}}$$

$$F^*(f)(d) = \hom(f, m(-,d)) \qquad F_*(g)(c) = \hom(m(c,-), g)$$

Category 2 instead of Set

Profunctor
(function) $r: X \times Y \longrightarrow \{0, 1\}$

$$X \xrightarrow{\ r\ } \{0,1\}^Y \qquad R^*$$
$$X \longrightarrow \{0,1\}^X$$

$$0 \rightarrow 1$$
$$\boxed{\textcircled{A}\ B}\ A \rightsquigarrow B$$

$$R^*(A) = \{y \in Y : R(x, y) = 1 \text{ for all } x \in A\}$$

$$\text{Set}^Y \longleftrightarrow Y$$
$$R_* \downarrow \qquad \swarrow r$$
$$\text{Set}^X$$

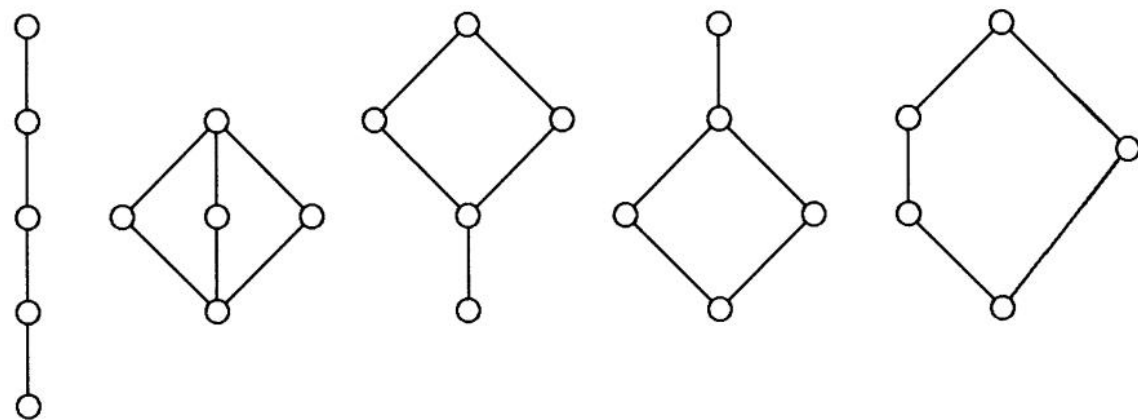$$R_*(B) = \{x \in X : R(x, y) = 1 \text{ for all } y \in B\}$$

Fixed obj of $R^* R_*$ and $R_* R^*$ are formal concept.

# Formal concept analysis

- 形式概念分析（Formal Concept Analysis，简称FCA）是一种基于格论（Lattice Theory）应用分支的数学理论，它使得概念的形式化成为可能，并以对象-属性的形式分析数据。FCA起源于20世纪80年代初，由Rudolf Wille提出，并在过去的几十年中逐渐发展成为一个知识表示和数据分析工具，广泛应用于信息检索、机器学习、数据挖掘、知识发现、文本挖掘等领域。

- FCA的一个关键优势是它提供了一种直观的方式来表示和分析复杂的数据结构，通过概念格图（Hasse Diagram）可以清晰地展示对象和属性之间的层次关系。此外，FCA还能够处理多值属性和不完全数据，使其成为一个灵活和强大的数据分析工具

**Definition 9.** Let $(M, \leq)$ be an ordered set and $A$ a subset of $M$. A **lower bound** of $A$ is an element $s$ of $M$ with $s \leq a$ for all $a \in A$. An **upper bound** of $A$ is defined dually. If there is a largest element in the set of all lower bounds of $A$, it is called the **infimum** of $A$ and is denoted by $\inf A$ or $\bigwedge A$; dually, a least upper bound is called **supremum** and denoted by $\sup A$ or $\bigvee A$. If $A = \{x, y\}$, we also write $x \wedge y$ for $\inf A$ and $x \vee y$ for $\sup A$. Infimum and supremum are frequently also called **meet** and **join**. $\diamond$

**Definition 10.** An ordered set $V := (V, \leq)$ is a **lattice**, if for any two elements $x$ and $y$ in $V$ the supremum $x \vee y$ and the infimum $x \wedge y$ always exist.

**Definition 18.** A **formal context** $\mathbb{K} := (G, M, I)$ consists of two sets $G$ and $M$ and a relation $I$ between $G$ and $M$. The elements of $G$ are called the **objects** and the elements of $M$ are called the **attributes** of the context[1]. In order to express that an object $g$ is in a relation $I$ with an attribute $m$, we write $gIm$ or $(g, m) \in I$ and read it as "the object $g$ **has** the attribute $m$".

|   |            | a | b | c | d | e | f | g | h | i |
|---|------------|---|---|---|---|---|---|---|---|---|
| 1 | Leech      | × | × |   |   |   |   | × |   |   |
| 2 | Bream      | × | × |   |   |   |   | × | × |   |
| 3 | Frog       | × | × | × |   |   |   | × | × |   |
| 4 | Dog        | × |   | × |   |   |   | × | × | × |
| 5 | Spike – weed | × | × |   | × |   | × |   |   |   |
| 6 | Reed       | × | × | × | × |   | × |   |   |   |
| 7 | Bean       | × |   | × | × | × |   |   |   |   |
| 8 | Maize      | × |   | × | × |   | × |   |   |   |

**Figure 1.1** Context of an educational film "Living Beings and Water". The attributes are: a: needs water to live, b: lives in water, c: lives on land, d: needs chlorophyll to produce food, e: two seed leaves, f: one seed leaf, g: can move around, h: has limbs, i: suckles its offspring.

**Definition 19.** For a set $A \subseteq G$ of objects we define

$$A' := \{m \in M \mid gIm \text{ for all } g \in A\}$$

(the set of attributes common to the objects in $A$). Correspondingly, for a set $B$ of attributes we define

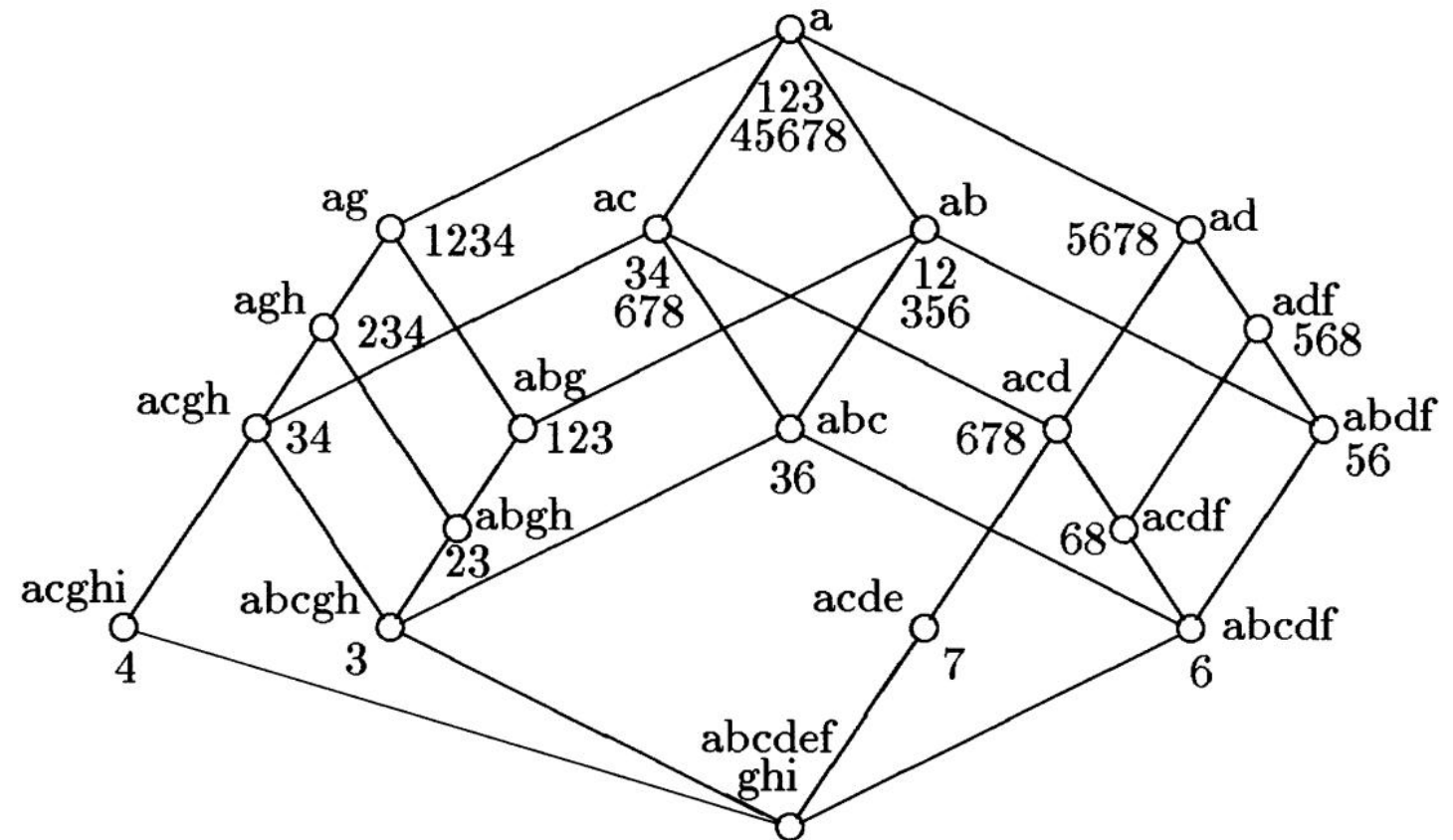$$B' := \{g \in G \mid gIm \text{ for all } m \in B\}$$

(the set of objects which have all attributes in $B$).[2] $\diamond$

**Definition 20.** A **formal concept** of the context $(G, M, I)$ is a pair $(A, B)$ with $A \subseteq G, B \subseteq M$, $A' = B$ and $B' = A$. We call $A$ the **extent** and $B$ the **intent** of the concept $(A, B)$. $\mathfrak{B}(G, M, I)$ denotes the set of all concepts of the context $(G, M, I)$. $\diamond$
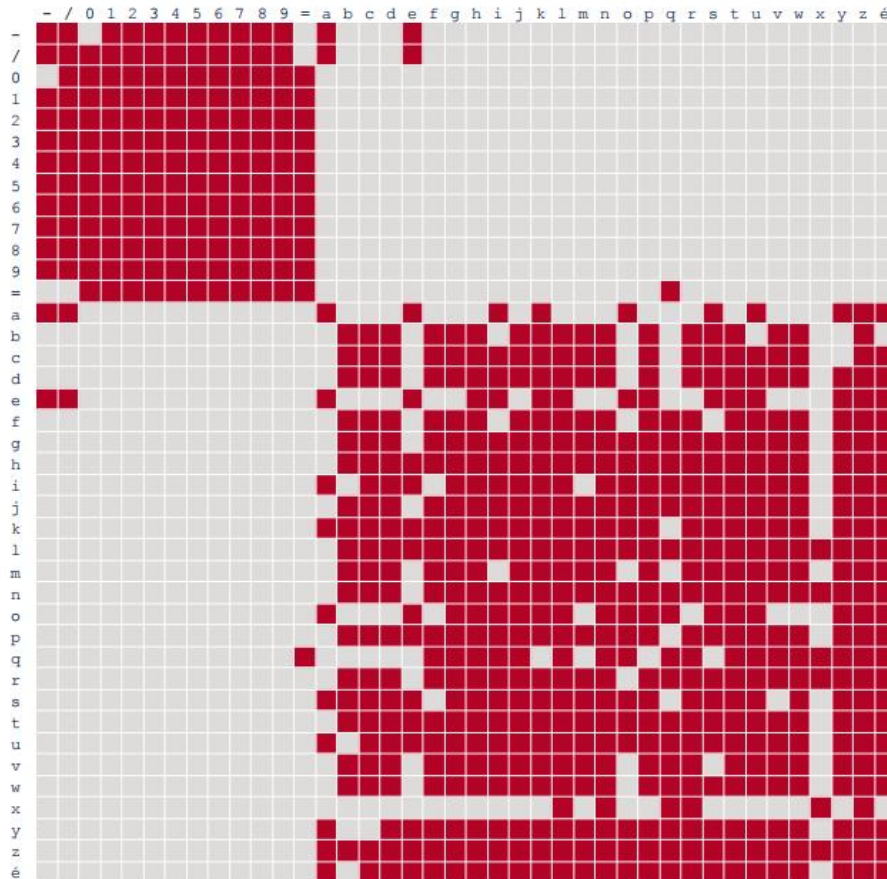
**Definition 21.** If $(A_1, B_1)$ and $(A_2, B_2)$ are concepts of a context, $(A_1, B_1)$ is called a **subconcept** of $(A_2, B_2)$, provided that $A_1 \subseteq A_2$ (which is equivalent to $B_2 \subseteq B_1$). In this case, $(A_2, B_2)$ is a **superconcept** of $(A_1, B_1)$, and we write $(A_1, B_1) \leq (A_2, B_2)$. The relation $\leq$ is called the **hierarchical order** (or simply **order**) of the concepts. The set of all concepts of $(G, M, I)$ ordered in this way is denoted by $\underline{\mathfrak{B}}(G, M, I)$ and is called the **concept lattice** of the context $(G, M, I)$.
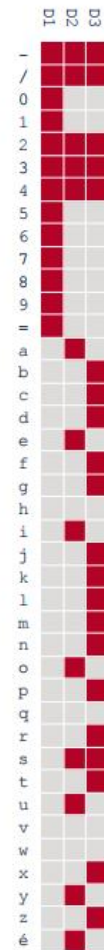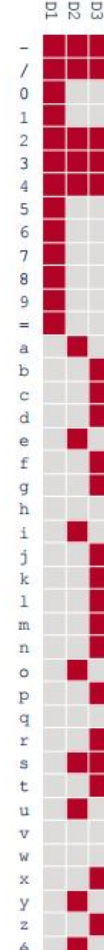
# Binary: Formal Concepts
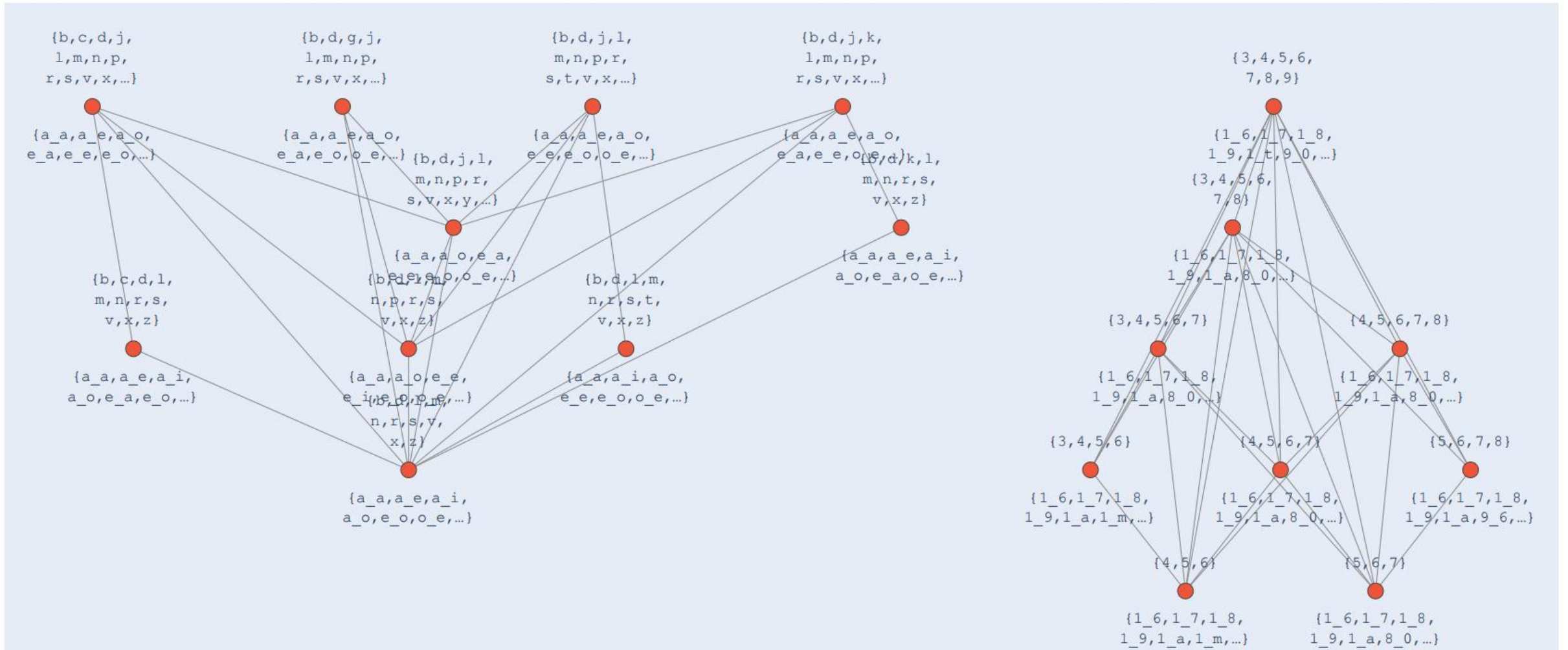
$$(M \times M^*)v = \lambda v$$

# Formal Concepts Words