

Corpus Analysis from a Mathematical Perspective

Corpus Statistics Research Group launch event
Birmingham, 11th Feb 2016

Simon Preston (University of Nottingham)

Joint work with R. Carrington, A. Hennessey, M. Mahlberg, K.
Severn, Y. van Gennip, V. Wiegand



Corpus as a mathematical object

THE POSTHUMOUS PAPERS
OF
THE PICKWICK CLUB

CHAPTER I
THE PICKWICKIANS

The first ray of light which illumines the gloom, and converts into a dazzling brilliancy that obscurity in which the earlier history of the public career of the immortal Pickwick would appear to be involved, is derived from the perusal of the following entry in the Transactions of the Pickwick Club, which the editor of these papers feels the highest pleasure in laying before his readers, as a proof of the careful attention, indefatigable assiduity, and nice discrimination, with which his search among the multifarious documents confided to him has been conducted.

'May 12, 1827. Joseph Smiggers, Esq., P.V.P.M.P.C. [Perpetual Vice-President--Member Pickwick Club], presiding. The following resolutions unanimously agreed to:--

'That this Association has heard read, with feelings of unmingled satisfaction, and unqualified approval, the paper communicated by Samuel Pickwick, Esq., G.C.M.P.C. [General Chairman--Member Pickwick Club], entitled "Speculations on the Source of the Hampstead Ponds, with some Observations on the Theory of Tittlebats;" and that this Association does hereby return its warmest thanks to the said Samuel Pickwick, Esq., G.C.M.P.C., for the same.

'That while this Association is deeply sensible of the advantages which must accrue to the cause of science, from the production to which they have just adverted--no less than from the unwearied researches of Samuel Pickwick, Esq., G.C.M.P.C., in Hornsey, Highgate, Brixton, and Camberwell--they cannot but entertain

OLIVER TWIST
OR
THE PARISH BOY'S PROGRESS
BY
CHARLES DICKENS

CHAPTER I

TREATS OF THE PLACE WHERE OLIVER TWIST WAS BORN
CIRCUMSTANCES ATTENDING HIS BIRTH

Among other public buildings in a certain town, for reasons it will be prudent to refrain from mentioning, which I will assign no fictitious name, there is a workhouse, common to most towns, great or small: to wit, in this workhouse was born; on a day and date I will not trouble myself to repeat, inasmuch as it can be of no consequence to the reader, in this stage of the events; the item of mortality whose name is printed at the end of this chapter.

For a long time after it was ushered into this world, and trouble, by the parish surgeon, it remained in considerable doubt whether the child would survive, and name at all; in which case it is somewhat more than probable these memoirs would never have appeared; or, in being comprised within a couple of pages, they would have possessed the inestimable merit of being the most faithful specimen of biography, extant in the literature of the age or country.

Although I am not disposed to maintain that the birth of Oliver Twist, in this workhouse, is in itself the most fortunate and extraordinary circumstance that can possibly befall a human being, I will say that in this particular instance, it was the only circumstance that could by possibility have occurred to Oliver Twist that could by possibility have occurred to him.

Corpus

THE POSTHUMOUS PAPERS
OF
THE FIDELITY CLUB

CHAPTER I
THE FIDELITY CLUB

The first ray of light which illumines the gloom, and converts
into a dazzling brilliancy that obscurity in which the earlier
history of the public career of the immortal Pickwick would
appear to be involved, is derived from the perusal of the following
entry in the Transactions of the Pickwick Club, which the editor
of these papers feels the highest pleasure in laying before his
readers, as a proof of the careful attention, indefatigable assiduity,
and nice discrimination, with which his search among the multivolume
documents confided to him has been conducted.

"May 12, 1827. Joseph Seagraves, Esq., F.R.S.E.P.C., President
Vice-President-Nether Pickwick Club, presiding. The following
resolutions unanimously agreed to:-

"That this Association has heard read, with feelings of unqualified
satisfaction, and unqualified approval, the paper communicated by Samuel
Pickwick, Esq., G.C.M.P.C., General Chairman-Nether Pickwick Club,
entitled "Speculations on the Source of the Hemphill Ponds, with some
Observations on the Theory of Titled Subjects;" and that this Association
does hereby return its warmest thanks to the said Samuel
Pickwick, Esq., G.C.M.P.C., for the same.

"That while this Association is deeply sensible of the advantages
which must accrue to the cause of science, from the production
to which they have just admitted-re less than from the unwarmed
responses of Samuel Pickwick, Esq., G.C.M.P.C., in memory,
Migajato, Brivato, and Camberwell-they cannot but entertain
the opinion, that the said paper is a masterpiece of science."

The old Gossamer Shop
By Charles Dickens

CHAPTER I

Night is generally my time for walking. In the summer I often leave
home early in the morning, and ram about Italia and Luna all day,
or even escape for days or weeks together; but, saving in the
country, I seldom go out until after dark, though, heaven be
thanked, I love its light and feel the cheerfulness it sheds upon the
earth, as much as any creature living.

I have fallen seemingly into this habit, both because it furnishes my
solace and because it affords me greater opportunity of speculating on
the characters and actions of those who fill the streets, the
glare and hurry of broad noon are not adapted to lofty pursuits like
mine; a glimpse of evening leaves caught by the light of a street lamp
or a lamp window is often better for me purpose than their full
revelation in the daylight, and, if I meet and often destroy an air-built castle
in this respect than day, which has often destroys an air-built castle
at the moment of its completion, without the least ceremony or remorse.

That constant pacing to and fro, that never-ending restlessness, that
incessant tread of feet meeting the rough stones smooth and glossy-is it
not a wonder how the dwellers in narrow ways can bear to sleep?
It is this of a sick man in such a place as Saint Martin's Court,
listening to the footstep, and in the midst of pain and weariness
obliged, despite himself, to watch it; it were a task he must perform
to detect the child's step from the man's, the slithered beggar from
the booted exquisites, the Gipsy from the Gipsy, the old deaf
of the snuffling ancient from the quick tread of an expectant
pensioner-nether-thing of the hum and noise always being present to his
sense-out of the stream of life, that will not allow himself to be
any one

SILVER TWIST
OR
THE PARISH BOY'S PROGRESS
BY
CHARLES DICKENS

CHAPTER I

TOWARDS OF THE PLACE WHERE OLIVER TWIST WAS BORN AND OF THE
CIRCUMSTANCES ATTENDING HIS BIRTH

Among other public buildings in a certain town, which for many
years I still remember to refrain from mentioning, and to
which I will assign no fictitious name, there is an anciently
common to most towns, great or small: is wit, a workhouse; and
in this workhouse was born, on a day and date which I need not
trouble myself to repeat, inasmuch as it can be of no possible
consequence to the reader, in this stage of the business as all
events; the time of mortality whose name is prefixed to the head
of this chapter.

For a long time after it was ushered into this world of sorrow
and trouble, by the parish surgeon, it remained a matter of
considerable doubt whether the child would survive to bear any
name at all; in which case it is somewhat more than probable that
these memoirs would never have appeared; or, if they had, that
being comprised within a couple of pages, they would have
possessed the unsatisfactory merit of being the most concise and
faithful specimen of biography, extant in the literature of any
age or country.

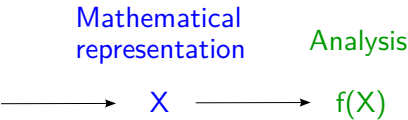
Although I am not disposed to maintain that the being born in a
workhouse, is in itself the most fortunate and enviable
circumstance that can possibly befall a human being, I do mean to
say that in this particular instance, it was the best thing for
Oliver Twist that could be sensibly have occurred. The fact

THE LIFE AND ADVENTURES OF NICHOLAS NICKLEBY
CONTAINING A FAITHFUL ACCOUNT OF THE PERIPATETIC, Misfortunes,
Grievings, Downfallings and Complete Career of the Nickleby Family
by Charles Dickens

AUTHOR'S PREFACE

"This story was begun, within a few months after the publication of
the completed "Pickwick Papers." There were, then, a good many cheap
furniture shops in existence. There are very few now.

Of the monstrous neglect of education in England, and the disregard
of it by the State as a means of forming good or bad citizens; and
misdeeds or happy men, private schools had offered a notable
example. Although any man who had proved his aptness for any other
pursuit in life, was free, without examination or qualification,
to open a school anywhere; although preparation for the functions he
undertook, was required in the surgeon who insisted on being a boy
into the world, or negro and day assistant, perhaps, to send him out of
it in the evening, the attorney, the butcher, the baker, the
confectioner maker; the whole round of crafts and trades, the
schoolmaster existed; and although schoolmasters, as a race, were
the bloodiest and ugliest who might naturally be expected to
spring from such a state of things, and to flourish in it; these
fortunate schoolmasters were the laziest and most rotten round in the
whole ladder. Trade in the market, indifference, or indelicacy of
parents, and the helplessness of children, however, were, I should
say, to show few considerate persons would have extracted the hard
and tedious of a horse at a dog; they formed the worthy core of
a structure, which, for abundance and a magnificent high-minded
LORDS-ALAN neglect, has rarely been exceeded in the world."



Corpus analysis

Corpus

THE POSTHUMOUS PAPERS
OF
THE FIDELITY CLUB

CHAPTER I
THE FIDELITY CLUB

The first ray of light which illumines the gloom, and converts into a dazzling brilliancy that obscurity in which the warlike history of the public career of the immortal Pickwick would appear to be involved, is derived from the perusal of the following entry in the Transactions of the Fidelity Club, which the editor of these papers feels the highest pleasure in laying before his readers, as a proof of the careful attention, indefatigable assiduity, and nice discrimination, with which his search among the multivolume documents confided to him has been conducted.

May 12, 1827. Joseph Snuggles, Esq., F.R.S.N.P.C. (Presidential Vice-President-Nether Pickwick Club), presiding. The following resolutions unanimously agreed to:-

"That this Association has heard read, with feelings of unqualified satisfaction, and unqualified approval, the paper communicated by Samuel Pickwick, Esq., G.C.M.P.C. (General Chairman-Nether Pickwick Club), entitled 'Speculations on the Source of the Hemphill Ponds, with some Observations on the Theory of Titled Subjects'; and that this Association does hereby return its warmest thanks to the said Samuel Pickwick, Esq., G.C.M.P.C., for the same.

"That while this Association is deeply sensible of the advantages which must accrue to the cause of science, from the production to which they have just admitted-re less than from the unwarmed recollections of Samuel Pickwick, Esq., G.C.M.P.C. in memory, mythology, fiction, and cabalwell-they cannot but entertain

SILVER TWIST
OR
THE PARISH BOY'S PROGRESS
BY
CHARLES DICKENS

CHAPTER I

YEARS OF THE PLACE WHERE OLIVER TWIST WAS BORN AND OF THE CIRCUMSTANCES ATTENDING HIS BIRTH

Among other public buildings in a certain town, which for many reasons it will be prudent to refrain from mentioning, and to which I will assign no fictitious name, there is an anciently common to most towns, great or small: it is, a workhouse; and in this workhouse was born, on a day and date which I need not trouble myself to repeat, inasmuch as it can be of no possible consequence to the reader, in this stage of the business as at all events; the time of mortality whose name is prefixed to the head of this chapter.

For a long time after it was ushered into this world of sorrow and trouble, by the parish surgeon, it remained a matter of considerable doubt whether the child would survive to bear any name at all; in which case it is somewhat more than probable that these memoirs would never have appeared; or, if they had, that being comprised within a couple of pages, they would have possessed the unsatisfactory world of being the most concise and faithful specimen of biography, extant in the literature of any age or country.

Although I am not disposed to maintain that the being born in a workhouse, is in itself the most fortunate and auspicious circumstance that can possibly befall a human being, I do mean to say that in this particular instance, it was the best thing for Oliver Twist that could be sensibly be have occurred. The fact

The Old Curiosity Shop

By Charles Dickens

CHAPTER I

Night is generally my time for walking. In the summer I often leave home early in the morning, and ram about Italia and Luna all day, or even escape for days or weeks together; but, saving in the country, I seldom go out until after dark, though, when he is thence, I love its lights and feel the cheerfulness it sheds upon the earth, as much as any creature living.

I have fallen seemingly into this habit, both because it favours my idleness and because it affords me greater opportunity of speculating on the characters and accidents of those who fill the streets, the glare and hurry of broad noon are not adapted to jobs peripatetic like mine; a glimpse of passing faces caught by the light of a street lamp, or a single window is often better for my purpose than their full revelation in the daylight, and, if I meet and stare, I look in kinder in this respect than day, which has often destroyed an air-built castle at the moment of its completion, without the least offence or remorse.

This constant peering to and fro, this never-ending restlessness, that incessant thread of feet meeting the rough stones smooth and glossy-is it not a wonder how the dwellers in narrow ways can bear to live? Think of a sick man in such a place as Saint Martin's Court, listening to the footstep, and in the midst of pain and weariness obliged, despite himself, to think it were a task he must perform to detect the child's step from the man's, the slighted beggar from the bearded exquisites, the foreigner from the Brit, the dark devil of the snarling adrover from the quick tread of an expectant vice-president-Nether Pickwick Club! The hum and noise always being present to his sense, and of the stream of life, that will not allow, because on an

THE LIFE AND ADVENTURES OF NICHOLAS NICKLEBY, CONTAINING A FAITHFUL ACCOUNT OF THE PERIPATICS, Misfortunes, Extraneous, Downfallings and Complete Career of the Nickleby Family

by Charles Dickens

AUTHOR'S PREFACE

This story was begun, within a few months after the publication of the completed 'Pickwick papers.' There were, then, a good many cheap

of the monstrous neglect of education in England, and the disregard of it by the state as a means of forming good or bad citizens; and miserable or happy men, private schools long offered a notable example. Although any man who had proved his aptitude for any other profession in life, was free, without examination or qualification, to open a school anywhere; although preparation for the functions he undertook, was not made in the manner he destined to bring a boy into the world, the children, the scholars, the laborer, the collector's man; the whole round of crafts and trades, the schoolmaster existed; and although schoolmasters, as a race, were the bloodiest and ugliest who might naturally be expected to spring from such a state of things, and to flourish in it; these fortunate schoolmasters were the kindest and most rather round in the whole ladder. Traders in the market, and therefore, or inability of parents, and the helplessness of children, however, were, to show few considerate persons would have entrusted the board and lodgings of a horse at a shop; they formed the worthy cornerstone of a structure, which, for abundance and magnificent high-minded LONDON-ALAN neglect, has rarely been exceeded in the world.

Mathematical
representation

Analysis



Analysis = studying patterns
- checking one is really there
- identifying new ones

Why is this perspective helpful?

- Deciding on X forces us to decide:
 - what in the corpus is important
 - what we are happy to discard
- For a given X we have a “toolbox” of available methods from which to choose $f(X)$: the **abstraction is powerful**.
- It helps us understand the $f(X)$ we choose to use.
- ... which is essential for developing new methodologies.

Example: Dickens novels

X as “bag of words” representation.

| | said | one | will | now | little | poor | upon | mrs | |
|------------|------|------|------|-----|--------|------|------|------|-----|
| <i>PP</i> | 3321 | 766 | 437 | 471 | 651 | 95 | 608 | 508 | |
| <i>OT</i> | 1232 | 457 | 302 | 280 | 276 | 97 | 477 | 264 | |
| <i>NN</i> | 2706 | 1019 | 712 | 608 | 743 | 262 | 1065 | 1040 | ... |
| <i>OCS</i> | 1420 | 653 | 331 | 436 | 646 | 177 | 796 | 252 | |
| <i>BR</i> | 1454 | 839 | 401 | 509 | 391 | 136 | 911 | 189 | |
| <i>MC</i> | 2786 | 1042 | 629 | 705 | 686 | 150 | 1153 | 953 | |
| <i>DS</i> | 2561 | 921 | 578 | 713 | 943 | 199 | 1105 | 1333 | |
| <i>DC</i> | 2950 | 908 | 531 | 741 | 1096 | 187 | 806 | 673 | ... |
| <i>BH</i> | 1743 | 971 | 805 | 909 | 1152 | 230 | 786 | 677 | |
| <i>HT</i> | 727 | 292 | 233 | 268 | 200 | 58 | 285 | 392 | |
| <i>LD</i> | 2139 | 1000 | 663 | 661 | 1454 | 261 | 779 | 928 | |
| <i>TTC</i> | 661 | 438 | 290 | 262 | 267 | 87 | 289 | 18 | |
| <i>GE</i> | 1349 | 502 | 174 | 453 | 371 | 77 | 366 | 164 | ... |
| <i>OMF</i> | 2180 | 859 | 622 | 757 | 878 | 252 | 753 | 988 | |
| <i>MED</i> | 406 | 229 | 266 | 206 | 203 | 70 | 227 | 77 | |

Such a “data matrix” is the central object in statistical
multivariate analysis.

Analysis method: matrix factorisation

- Break down X into the product “A times B”:

$$\begin{matrix} \text{novel} \times \text{word} \\ X \end{matrix} \approx \begin{matrix} \text{novel} \times r \\ A \end{matrix} \times \begin{matrix} r \times \text{word} \\ B \end{matrix}$$

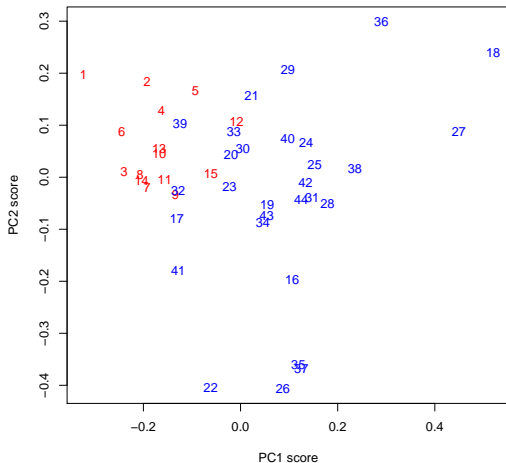
Analysis method: matrix factorisation

- Break down X into the product “A times B”:

$$\begin{matrix} \text{novel} \times \text{word} \\ X \end{matrix} \approx \begin{matrix} \text{novel} \times r \\ A \end{matrix} \times \begin{matrix} r \times \text{word} \\ B \end{matrix}$$

- Rows of B represent “features” found in corpus.
- Rows of A represent novels as “scores” for these features.
- Different constraints on A and B results in well-known methods:
 - Principal component analysis (PCA)
 - Latent semantic analysis
 - Non-negative matrix factorisation
 - (Topic modelling)

PCA for Dickens and other 19C novels



Red = Dickens novels (numbering indicates chronology)

Blue = Misc other 19C novels (numbering arbitrary)

PC interpretation

- Interpretation of scores in A?
- First and second rows/features of B:

| Row 1 | | Row 2 | |
|-------|--------|--------|--------|
| said | -0.559 | miss | -0.424 |
| mrs | -0.184 | mrs | -0.274 |
| sir | -0.175 | much | -0.129 |
| old | -0.131 | must | -0.127 |
| upon | -0.125 | little | -0.112 |
| ⋮ | ⋮ | ⋮ | ⋮ |
| yet | 0.128 | man | 0.122 |
| will | 0.143 | upon | 0.193 |
| now | 0.146 | said | 0.256 |

Other representations?

"Citizen Evremonde," she said, touching him with her cold hand. "I am a **poor little** seamstress, who was with you in La Force."

He murmured for answer: "True. I forget what you were accused of?"

"Plots. Though the just Heaven knows that I am innocent of any. Is it likely? Who would think of plotting with a **poor little** weak creature like me?"

The forlorn smile with which she said it, so touched him, that tears started from his eyes.

"I am not afraid to die, Citizen Evremonde, but I have done nothing. I am not unwilling to die, if the Republic which is to do so much good to us **poor**, will profit by my death; but I do not know how that can be, Citizen Evremonde. Such a **poor weak little** creature!"

As the last thing on earth that his heart was to warm and soften to, it warmed and softened to this pitiable girl.

"I heard you were released, Citizen Evremonde. I hoped it was true?"

"It was. But, I was again taken and condemned."

"If I may ride with you, Citizen Evremonde, will you let me hold your hand? I am not afraid, but I am **little and weak**, and it will give me more courage."

(A Tale of Two Cities, Dickens)

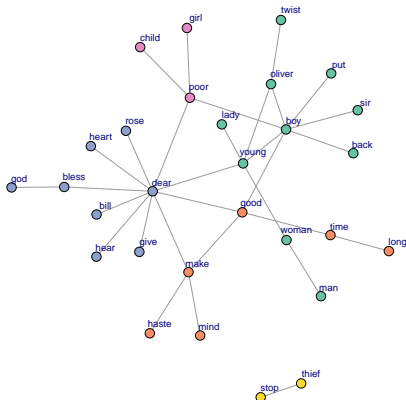
Speech from Oliver Twist: co-occurrence matrix

word \times word
X =

| | dear | boy | good | bill | hear | sir | give | lady | haste | girl | bless | mind | oliver | stop | young | back | make | child | long | man | woman | time | heart | poor | god | put | twist | rose | thief |
|--------|------|-----|------|------|------|-----|------|------|-------|------|-------|------|--------|------|-------|------|------|-------|------|-----|-------|------|-------|------|-----|-----|-------|------|-------|
| dear | 22 | 5 | 8 | 9 | 9 | 4 | 7 | 5 | 0 | 4 | 8 | 3 | 3 | 3 | 10 | 3 | 7 | 6 | 2 | 2 | 1 | 1 | 10 | 7 | 5 | 0 | 1 | 7 | 0 |
| boy | 5 | 20 | 10 | 1 | 2 | 8 | 3 | 0 | 0 | 2 | 0 | 5 | 9 | 0 | 7 | 7 | 2 | 3 | 1 | 4 | 1 | 0 | 0 | 10 | 0 | 8 | 3 | 0 | 0 |
| good | 8 | 10 | 16 | 1 | 6 | 3 | 2 | 0 | 0 | 1 | 2 | 2 | 0 | 0 | 6 | 1 | 7 | 1 | 2 | 2 | 1 | 11 | 0 | 2 | 0 | 0 | 1 | 1 | 0 |
| bill | 9 | 1 | 1 | 12 | 1 | 0 | 3 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 3 | 1 | 2 | 2 | 0 | 0 | 0 | 0 |
| hear | 9 | 2 | 6 | 1 | 12 | 1 | 1 | 2 | 0 | 0 | 2 | 4 | 0 | 0 | 2 | 0 | 3 | 0 | 1 | 3 | 2 | 1 | 1 | 2 | 0 | 1 | 0 | 1 | 0 |
| sir | 4 | 8 | 3 | 0 | 1 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 1 | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 2 | 1 | 3 | 0 | 0 | 2 | 0 | 0 |
| give | 7 | 3 | 2 | 3 | 1 | 0 | 10 | 3 | 1 | 1 | 0 | 2 | 2 | 1 | 2 | 3 | 2 | 0 | 2 | 3 | 0 | 3 | 1 | 1 | 1 | 2 | 0 | 0 | 0 |
| lady | 5 | 0 | 0 | 0 | 2 | 0 | 3 | 2 | 0 | 1 | 3 | 1 | 1 | 1 | 10 | 2 | 2 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 0 | 0 | 0 |
| haste | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| girl | 4 | 2 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 2 | 0 | 4 | 2 | 0 | 2 | 3 | 2 | 2 | 3 | 0 | 0 | 0 | 2 | 9 | 1 | 2 | 0 | 1 | 0 |
| bless | 8 | 0 | 2 | 0 | 2 | 0 | 0 | 3 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 6 | 2 | 8 | 0 | 0 | 0 | 0 |
| mind | 3 | 5 | 2 | 2 | 4 | 0 | 2 | 1 | 2 | 4 | 1 | 8 | 0 | 0 | 2 | 1 | 7 | 1 | 1 | 0 | 0 | 1 | 2 | 0 | 3 | 2 | 0 | 0 | 0 |
| oliver | 3 | 9 | 0 | 0 | 0 | 4 | 2 | 1 | 0 | 2 | 0 | 0 | 8 | 0 | 8 | 3 | 4 | 2 | 0 | 1 | 0 | 3 | 0 | 2 | 0 | 2 | 8 | 0 | 1 |
| stop | 3 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 8 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 7 |
| young | 10 | 7 | 6 | 1 | 2 | 1 | 2 | 10 | 0 | 2 | 0 | 2 | 8 | 0 | 8 | 0 | 2 | 2 | 1 | 5 | 7 | 5 | 4 | 6 | 0 | 4 | 4 | 0 | 0 |
| back | 3 | 7 | 1 | 0 | 0 | 0 | 3 | 2 | 0 | 3 | 1 | 1 | 3 | 0 | 0 | 6 | 3 | 0 | 4 | 6 | 0 | 5 | 1 | 2 | 3 | 1 | 0 | 0 | 2 |
| make | 7 | 2 | 7 | 0 | 3 | 0 | 2 | 2 | 9 | 2 | 0 | 7 | 4 | 0 | 2 | 3 | 4 | 3 | 2 | 3 | 2 | 3 | 1 | 0 | 0 | 2 | 0 | 1 | 1 |
| child | 6 | 3 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 2 | 0 | 1 | 2 | 0 | 2 | 0 | 3 | 4 | 1 | 4 | 0 | 0 | 2 | 7 | 1 | 1 | 1 | 1 | 0 |
| long | 2 | 1 | 2 | 1 | 1 | 1 | 2 | 0 | 0 | 3 | 1 | 1 | 0 | 1 | 1 | 4 | 2 | 1 | 4 | 1 | 0 | 7 | 2 | 1 | 1 | 0 | 0 | 0 | 1 |
| man | 2 | 4 | 2 | 1 | 3 | 1 | 3 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 5 | 6 | 3 | 4 | 1 | 2 | 7 | 3 | 3 | 4 | 0 | 1 | 0 | 0 | 1 |
| woman | 1 | 1 | 1 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 2 | 0 | 0 | 7 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| time | 1 | 0 | 11 | 3 | 1 | 2 | 3 | 1 | 0 | 0 | 0 | 1 | 3 | 0 | 5 | 5 | 3 | 0 | 7 | 3 | 0 | 6 | 2 | 2 | 0 | 5 | 0 | 0 | 1 |
| heart | 10 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 2 | 6 | 2 | 0 | 0 | 4 | 1 | 1 | 2 | 2 | 3 | 2 | 2 | 2 | 1 | 0 | 2 | 0 | 4 | 0 |
| poor | 7 | 10 | 2 | 2 | 2 | 3 | 1 | 0 | 0 | 9 | 2 | 0 | 2 | 0 | 6 | 2 | 0 | 7 | 1 | 4 | 0 | 2 | 1 | 2 | 0 | 2 | 1 | 0 | 0 |
| god | 5 | 0 | 0 | 2 | 0 | 0 | 1 | 4 | 0 | 1 | 8 | 3 | 0 | 1 | 0 | 3 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| put | 0 | 8 | 0 | 0 | 1 | 0 | 2 | 0 | 0 | 2 | 0 | 2 | 2 | 0 | 4 | 1 | 2 | 1 | 0 | 1 | 0 | 5 | 2 | 2 | 0 | 4 | 0 | 0 | 0 |
| twist | 1 | 3 | 1 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 0 | 4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| rose | 7 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 4 | 0 |
| thief | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 7 | 0 | 2 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 6 |

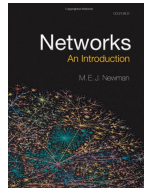
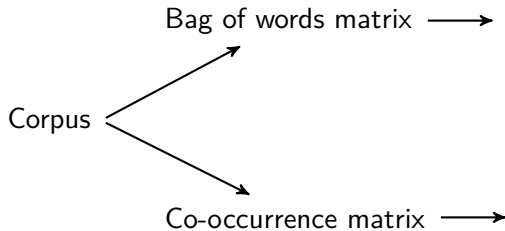
Speech from Oliver Twist: network visualisation

- Such matrix $\text{word} \times \text{word}$ X can be identified with a “graph” (network).
- Lots of methods available for graphs.
- → Yves’ talk later



Mathematical
representation, X

Analysis, $f(X)$



Challenges and directions

- How to analyse time structured corpora? (E.g. newspaper archive)
 - Bag of words approach: each row of X is associated with a time t_i , then consider time-weighted $X(t) \rightarrow$ Anthony's talk.
- How to harness tools of network analysis to analyse co-occurrence networks, e.g. clustering? \rightarrow Yves' talk.
- How to study time dependent networks? \rightarrow ongoing work.

Summary

- All methods of corpus analysis are a function $f(X)$ of a mathematical representation, X , of the corpus.
- Identifying X explicitly is **helpful**
 - to understand what information is used and what is discarded,
 - because abstraction provides a toolbox of methodologies, $f(X)$,
- ... and **essential**
 - to perform calculations for $f(X)$ efficiently,
 - **to develop new methodology**, extending existing $f(X)$.

Summary

- All methods of corpus analysis are a function $f(X)$ of a mathematical representation, X , of the corpus.
- Identifying X explicitly is **helpful**
 - to understand what information is used and what is discarded,
 - because abstraction provides a toolbox of methodologies, $f(X)$,
- ... and **essential**
 - to perform calculations for $f(X)$ efficiently,
 - **to develop new methodology**, extending existing $f(X)$.
- Many promising directions ahead!