

Terminology

Terminology

Unsupervised techniques:

Terminology

Unsupervised techniques: learning
(hidden or latent) structure in
unlabeled data.

e.g. PCA of legislators's votes:

Terminology

Unsupervised techniques: learning
(hidden or latent) structure in
unlabeled data.

e.g. PCA of legislators's votes: want to see
how they are organized—

Terminology

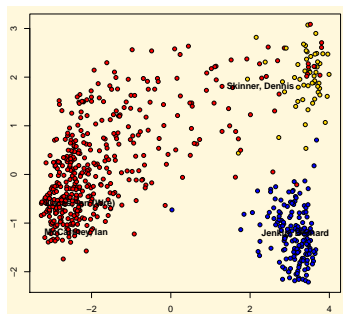
Unsupervised techniques: learning
(hidden or latent) structure in
unlabeled data.

e.g. PCA of legislators's votes: want to see
how they are organized—by party? by
ideology? by race?

Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

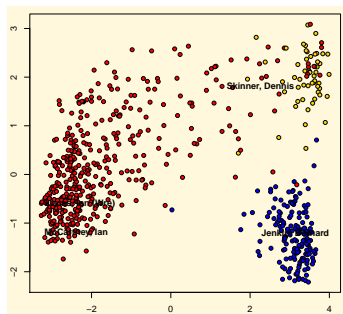
e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?

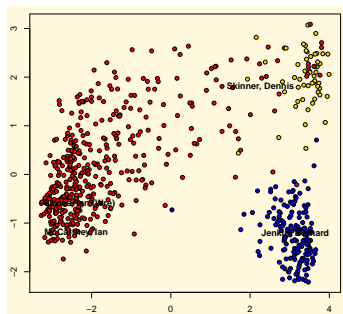


Supervised techniques:

Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?

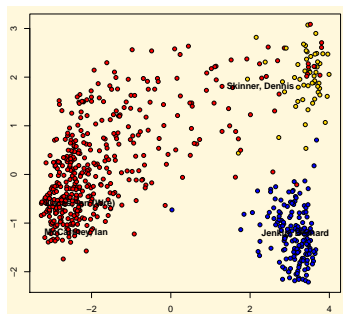


Supervised techniques: learning relationship between inputs and a labeled set of outputs.

Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



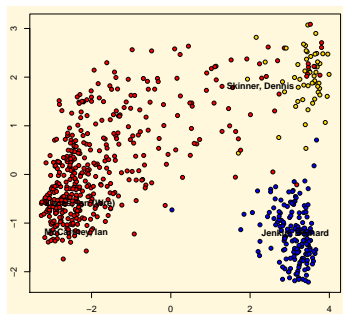
Supervised techniques: learning relationship between inputs and a labeled set of outputs.

e.g. opinion mining:

Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?



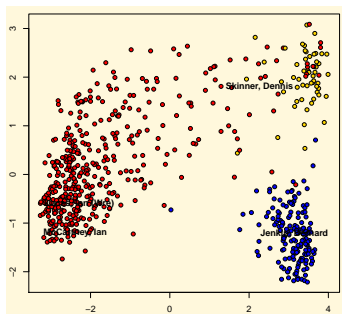
Supervised techniques: learning relationship between inputs and a labeled set of outputs.

e.g. opinion mining: what makes a critic like or dislike a movie ($y \in \{0, 1\}$)?

Terminology

Unsupervised techniques: learning (hidden or latent) structure in unlabeled data.

e.g. PCA of legislators's votes: want to see how they are organized—by party? by ideology? by race?




Supervised techniques: learning relationship between inputs and a labeled set of outputs.


e.g. opinion mining: what makes a critic like or dislike a movie ($y \in \{0, 1\}$)?


CRITIC REVIEWS FOR STAR WARS: EPISODE VII - THE FORCE AWAKENS

All Critics (313) | Top Critics (48) | My Critics | Fresh (293) | Rotten (20)


 The new movie, as an act of pure storytelling, streams by with fluency and zip.


[Full Review...](#) | December 21, 2015

 **Anthony Lane**
New Yorker
★ Top Critic


 At the end The Force Awakens looks more like a nostalgic film that will work as a transition to the new Star Wars' age. [Full Review in Spanish]


[Full Review...](#) | December 29, 2015

 **Salvador Franco Reyes**

 While Star Wars: The Force Awakens gets temporarily bogged down taking us back to the world that we left in 1983, it introduces us to the new and exciting torch-bearers of the franchise.

[Full Review...](#) | December 30, 2015

 **Blake Howard**
Graffiti With Punctuation

 This film is a well-planned product that balances nostalgia with the capacity to attract new generations into the Star Wars universe. [Full Review in Spanish]

[Full Review...](#) | December 29, 2015

Overview: Supervised Learning

Overview: Supervised Learning

label some examples of each category

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$)

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal,

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

train a 'machine' on these examples (e.g. logistic regression),

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

train a 'machine' on these examples (e.g. logistic regression), using the
features (DTM, other stuff) as the 'independent' variables.

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

train a 'machine' on these examples (e.g. logistic regression), using the **features** (DTM, other stuff) as the 'independent' variables.

e.g. does the commentator use the word 'fetus' or 'baby' in discussing abortion law?

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

train a 'machine' on these examples (e.g. logistic regression), using the **features** (DTM, other stuff) as the 'independent' variables.

e.g. does the commentator use the word 'fetus' or 'baby' in discussing abortion law?

classify use the learned relationship—some $f(x)$ —to predict the outcomes of documents ($y \in \{0, 1\}$, review sentiment)

Overview: Supervised Learning

label some examples of each category

e.g. some reviews that were positive ($y = 1$), some that were negative ($y = 0$);
some statements that were liberal, some that were conservative.

train a 'machine' on these examples (e.g. logistic regression), using the **features** (DTM, other stuff) as the 'independent' variables.

e.g. does the commentator use the word 'fetus' or 'baby' in discussing abortion law?

classify use the learned relationship—some $f(x)$ —to predict the outcomes of documents ($y \in \{0, 1\}$, review sentiment) not in the training set.

Overview: Dictionary

Overview: Dictionary

idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents

Overview: Dictionary

idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.

Overview: Dictionary

- idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.
- common in opinion mining/sentiment analysis,

Overview: Dictionary

- idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.
- common in opinion mining/sentiment analysis, and in coding events or manifestos.

Overview: Dictionary

- idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.
- common in opinion mining/sentiment analysis, and in coding events or manifestos.

Often **derived from** supervised learning techniques

Overview: Dictionary

idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.

- common in opinion mining/sentiment analysis, and in coding events or manifestos.

Often **derived from** supervised learning techniques
and often **used in** supervised learning problems, as a starting point.

Overview: Dictionary

idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.

- common in opinion mining/sentiment analysis, and in coding events or manifestos.

Often **derived from** supervised learning techniques
and often **used in** supervised learning problems, as a starting point.
so we'll cover them here.

Overview: Dictionary

idea: set of **pre-defined words** with specific connotations that allow us to **classify** documents automatically, quickly and accurately.

- common in opinion mining/sentiment analysis, and in coding events or manifestos.

Often **derived from** supervised learning techniques
and often **used in** supervised learning problems, as a starting point.
so we'll cover them here.

Classification with Dictionary Methods

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive',

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive', this speech is 'liberal'

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive', this speech is 'liberal'

- 2 Measure **extent to which** document is associated with given category

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive', this speech is 'liberal'

- 2 Measure **extent to which** document is associated with given category

e.g. this review is generally 'positive', but has some negative elements.

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive', this speech is 'liberal'

- 2 Measure **extent to which** document is associated with given category

e.g. this review is generally 'positive', but has some negative elements.

We have a **pre-determined** list of words, the (weighted) presence of which helps us with (1) and (2).

Classification with Dictionary Methods

Aim Typically we are trying to do one of two closely related things:

- 1 Categorize documents as belonging to a certain class (mutually exclusive? jointly exhaustive?)

e.g. this review is 'positive', this speech is 'liberal'

- 2 Measure **extent to which** document is associated with given category

e.g. this review is generally 'positive', but has some negative elements.

We have a **pre-determined** list of words, the (weighted) presence of which helps us with (1) and (2).

More Specifically

More Specifically

We have a set of **key words**, with attendant scores,

More Specifically

We have a set of **key words**, with attendant scores,
e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

More Specifically

We have a set of **key words**, with attendant scores,

- e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$
 - the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

$$\text{tone of document } i = \sum_{m=1}^M \frac{s_m w_{im}}{N_i}$$

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

$$\text{tone of document } i = \sum_{m=1}^M \frac{s_m w_{im}}{N_i}$$

where s_m is the score of word m

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

$$\text{tone of document } i = \sum_{m=1}^M \frac{s_m w_{im}}{N_i}$$

where s_m is the score of word m

and w_{im} is the number of occurrences of the m th dictionary word in the document i

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

$$\text{tone of document } i = \sum_{m=1}^M \frac{s_m w_{im}}{N_i}$$

where s_m is the score of word m

and w_{im} is the number of occurrences of the m th dictionary word in the document i

and N_i is the total number of all dictionary words in the document.

More Specifically

We have a set of **key words**, with attendant scores,

e.g. for movie reviews: 'terrible' is scored as -1 ; 'fantastic' as $+1$

→ the **relative rate** of occurrence of these terms tells us about the overall **tone** or category that the document should be placed in.

i.e. for document i and words $m = 1, \dots, M$ in the dictionary,

$$\text{tone of document } i = \sum_{m=1}^M \frac{s_m w_{im}}{N_i}$$

where s_m is the score of word m

and w_{im} is the number of occurrences of the m th dictionary word in the document i

and N_i is the total number of all dictionary words in the document.

→ just add up the number of times the words appear and multiply by the score (normalizing by doc dictionary presence)

(Simple) Example: Barnes' review of *The Big Short*

(Simple) Example: Barnes' review of *The Big Short*

Director and co-screenwriter Adam McKay (Step Brothers) bungles a great opportunity to savage the architects of the 2008 financial crisis in The Big Short, wasting an A-list ensemble cast in the process. Steve Carell, Brad Pitt, Christian Bale and Ryan Gosling play various tenuously related members of the finance industry, men who made made a killing by betting against the housing market, which at that point had superficially swelled to record highs. All of the elements are in place for a lacerating satire, but almost every aesthetic choice in the film is bad, from the U-Turn-era Oliver Stone visuals to Carell's sketch-comedy performance to the cheeky cutaways where Selena Gomez and Anthony Bourdain explain complex financial concepts. After a brutal opening half, it finally settles into a groove, and there's a queasy charge in watching a credit-drunk America walking towards that cliff's edge, but not enough to save the film.

Retain words in Hu & Liu Dictionary. . .

*Director and co-screenwriter Adam McKay (Step Brothers) bungles a **great** opportunity to **savage** the architects of the 2008 financial **crisis** in The Big Short, **wasting** an A-list ensemble cast in the process. Steve Carell, Brad Pitt, Christian Bale and Ryan Gosling play various **tenuously** related members of the finance industry, men who made made a **killing** by betting against the housing market, which at that point had **superficially swelled** to record highs. All of the elements are in place for a lacerating satire, but almost every aesthetic choice in the film is **bad**, from the U-Turn-era Oliver Stone visuals to Carell's sketch-comedy performance to the cheeky cutaways where Selena Gomez and Anthony Bourdain explain **complex** financial concepts. After a **brutal** opening half, it finally settles into a groove, and there's a queasy charge in watching a credit-**drunk** America walking towards that cliff's edge, but not **enough** to save the film.*

Retain words in Hu & Liu Dictionary. . .

great
crisis

savage
wasting

tenuously

killing

superficially swelled

bad

brutal

complex

drunk
enough

Simple math...

Simple math...

negative 11

Simple math...

negative 11

positive 2

Simple math...

negative 11

positive 2

total 13

Simple math...

negative 11

positive 2

total 13

$$\text{tone} = \frac{2-11}{13} = \frac{-9}{13}$$

Simple math...

negative 11

positive 2

total 13

$$\text{tone} = \frac{2-11}{13} = \frac{-9}{13}$$

