

Text as Data

Justin Grimmer

Professor
Department of Political Science
Stanford University

April 9th, 2019

Texts and Geometry

Consider a document-term matrix

$$\mathbf{x} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space**

Texts and Geometry

Consider a document-term matrix

$$\mathbf{x} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

- Provides a **geometry**

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

- Provides a **geometry** \rightsquigarrow modify with word weighting

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

- Provides a **geometry** \rightsquigarrow modify with word weighting
- Natural notions of **distance**

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

- Provides a **geometry** \rightsquigarrow modify with word weighting
- Natural notions of **distance**
- **Kernel Trick**: richer comparisons of large feature spaces

Texts and Geometry

Consider a document-term matrix

$$\mathbf{X} = \begin{pmatrix} 1 & 2 & 0 & \dots & 0 \\ 0 & 0 & 3 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & 0 & \dots & 3 \end{pmatrix}$$

Suppose documents live in a **space** \rightsquigarrow rich set of results from linear algebra

- Provides a **geometry** \rightsquigarrow modify with word weighting
- Natural notions of **distance**
- **Kernel Trick**: richer comparisons of large feature spaces
- Building block for clustering, supervised learning, and scaling

Texts in Space

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

Texts in Space

Doc1 = $(1, 1, 3, \dots, 5)$

Doc2 = $(2, 0, 0, \dots, 1)$

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

$$\text{Doc2} = (2, 0, 0, \dots, 1)$$

$$\mathbf{Doc1}, \mathbf{Doc2} \in \mathbb{R}^J$$

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

$$\text{Doc2} = (2, 0, 0, \dots, 1)$$

$$\mathbf{Doc1}, \mathbf{Doc2} \in \mathbb{R}^J$$

Inner Product between documents:

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

$$\text{Doc2} = (2, 0, 0, \dots, 1)$$

$$\mathbf{Doc1}, \mathbf{Doc2} \in \mathbb{R}^J$$

Inner Product between documents:

$$\mathbf{Doc1} \cdot \mathbf{Doc2} = (1, 1, 3, \dots, 5)' (2, 0, 0, \dots, 1)$$

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

$$\text{Doc2} = (2, 0, 0, \dots, 1)$$

$$\mathbf{Doc1}, \mathbf{Doc2} \in \mathbb{R}^J$$

Inner Product between documents:

$$\begin{aligned}\mathbf{Doc1} \cdot \mathbf{Doc2} &= (1, 1, 3, \dots, 5)' (2, 0, 0, \dots, 1) \\ &= 1 \times 2 + 1 \times 0 + 3 \times 0 + \dots + 5 \times 1\end{aligned}$$

Texts in Space

$$\text{Doc1} = (1, 1, 3, \dots, 5)$$

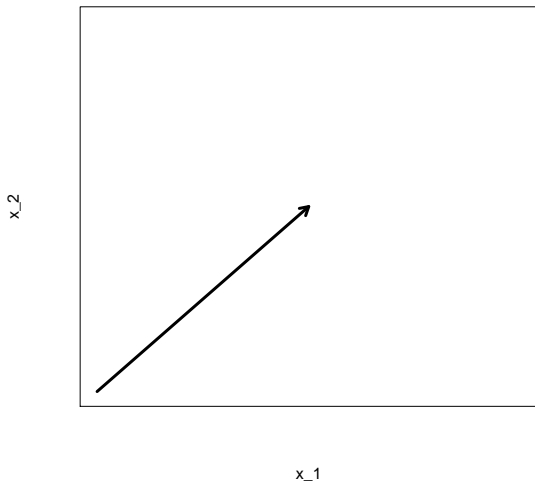
$$\text{Doc2} = (2, 0, 0, \dots, 1)$$

$$\mathbf{Doc1}, \mathbf{Doc2} \in \mathbb{R}^J$$

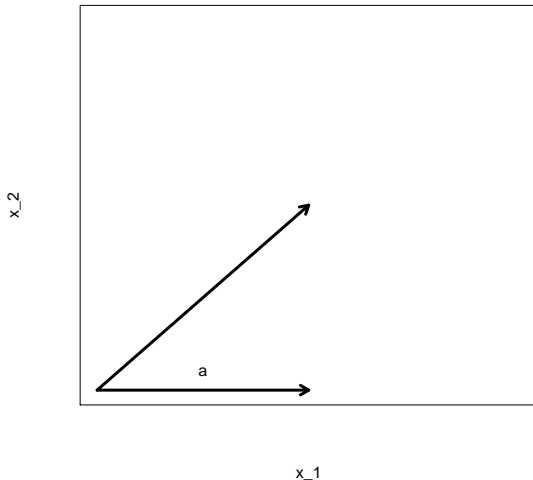
Inner Product between documents:

$$\begin{aligned}\mathbf{Doc1} \cdot \mathbf{Doc2} &= (1, 1, 3, \dots, 5)' (2, 0, 0, \dots, 1) \\ &= 1 \times 2 + 1 \times 0 + 3 \times 0 + \dots + 5 \times 1 \\ &= 7\end{aligned}$$

Vector Length

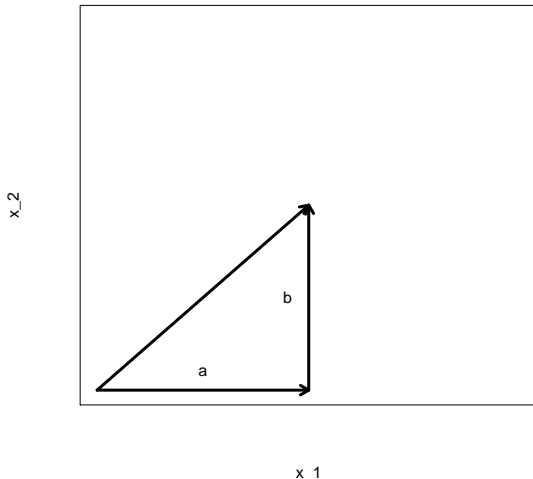


Vector Length



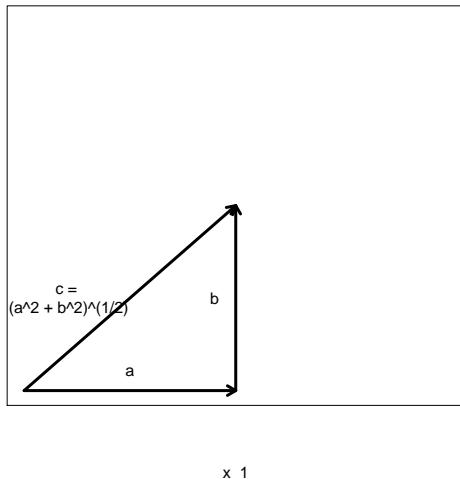
- **Pythagorean Theorem:**
Side with length a

Vector Length



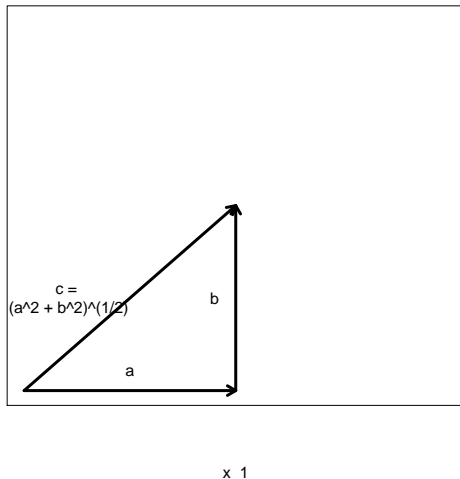
- **Pythagorean Theorem:**
Side with length a
- Side with length b and
right triangle

Vector Length



- **Pythagorean Theorem:**
Side with length a
- Side with length b and right triangle
- $c = \sqrt{a^2 + b^2}$

Vector Length



- **Pythagorean Theorem:**
Side with length a
- Side with length b and
right triangle
- $c = \sqrt{a^2 + b^2}$
- **This is generally true**

Vector (Euclidean) Length

Definition

Suppose $\mathbf{v} \in \mathbb{R}^J$. Then, we will define its *length* as

$$\begin{aligned}\|\mathbf{v}\| &= (\mathbf{v} \cdot \mathbf{v})^{1/2} \\ &= (v_1^2 + v_2^2 + v_3^2 + \dots + v_J^2)^{1/2}\end{aligned}$$

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

\rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

$$1) d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$$

\rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

1) $d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$

2) $d(\mathbf{X}_i, \mathbf{X}_j) = 0$ if and only if $\mathbf{X}_i = \mathbf{X}_j$

\rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

- 1) $d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$
- 2) $d(\mathbf{X}_i, \mathbf{X}_j) = 0$ if and only if $\mathbf{X}_i = \mathbf{X}_j$
- 3) $d(\mathbf{X}_i, \mathbf{X}_j) = d(\mathbf{X}_j, \mathbf{X}_i)$

\rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

- 1) $d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$
- 2) $d(\mathbf{X}_i, \mathbf{X}_j) = 0$ if and only if $\mathbf{X}_i = \mathbf{X}_j$
- 3) $d(\mathbf{X}_i, \mathbf{X}_j) = d(\mathbf{X}_j, \mathbf{X}_i)$
- 4) $d(\mathbf{X}_i, \mathbf{X}_k) \leq d(\mathbf{X}_i, \mathbf{X}_j) + d(\mathbf{X}_j, \mathbf{X}_k)$

\rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

- 1) $d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$
- 2) $d(\mathbf{X}_i, \mathbf{X}_j) = 0$ if and only if $\mathbf{X}_i = \mathbf{X}_j$
- 3) $d(\mathbf{X}_i, \mathbf{X}_j) = d(\mathbf{X}_j, \mathbf{X}_i)$
- 4) $d(\mathbf{X}_i, \mathbf{X}_k) \leq d(\mathbf{X}_i, \mathbf{X}_j) + d(\mathbf{X}_j, \mathbf{X}_k)$

Explore distance functions to compare documents \rightsquigarrow

Measures of Dissimilarity

Initial guess \rightsquigarrow Distance metrics

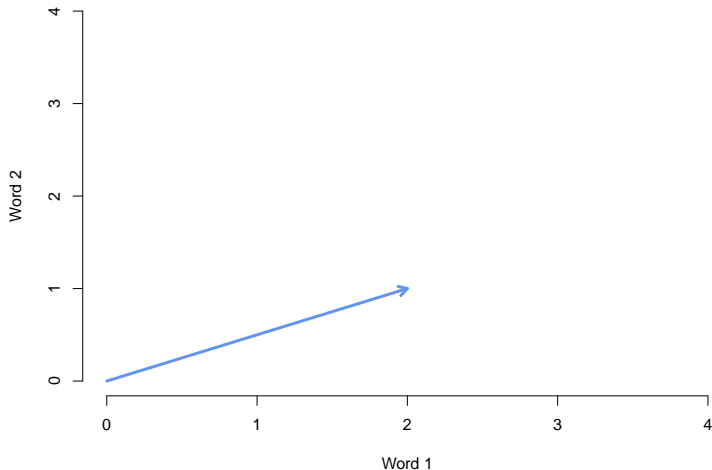
Properties of a metric: (distance function) $d(\cdot, \cdot)$. Consider arbitrary documents $\mathbf{X}_i, \mathbf{X}_j, \mathbf{X}_k$

- 1) $d(\mathbf{X}_i, \mathbf{X}_j) \geq 0$
- 2) $d(\mathbf{X}_i, \mathbf{X}_j) = 0$ if and only if $\mathbf{X}_i = \mathbf{X}_j$
- 3) $d(\mathbf{X}_i, \mathbf{X}_j) = d(\mathbf{X}_j, \mathbf{X}_i)$
- 4) $d(\mathbf{X}_i, \mathbf{X}_k) \leq d(\mathbf{X}_i, \mathbf{X}_j) + d(\mathbf{X}_j, \mathbf{X}_k)$

Explore distance functions to compare documents \rightsquigarrow Do we want additional assumptions/properties?

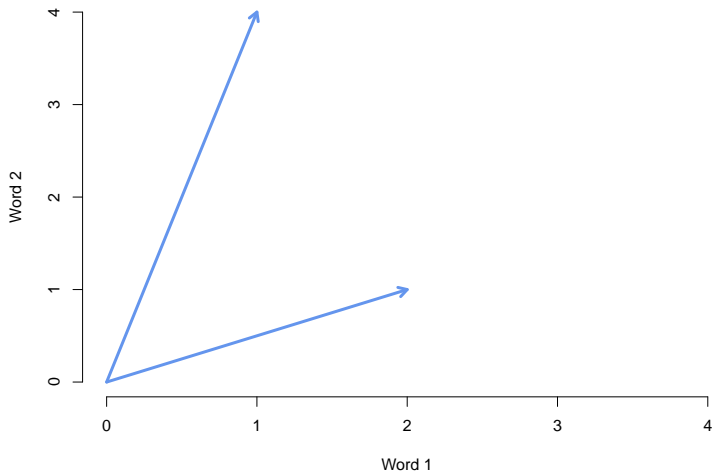
Measuring the Distance Between Documents

Euclidean Distance



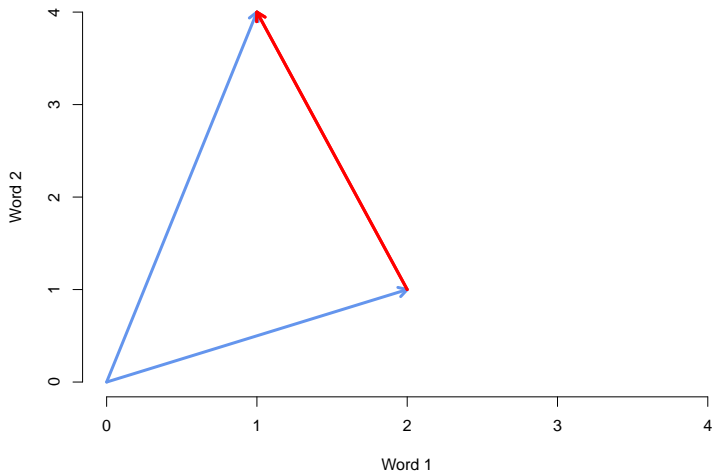
Measuring the Distance Between Documents

Euclidean Distance



Measuring the Distance Between Documents

Euclidean Distance



Measuring the Distance Between Documents

Definition

The Euclidean distance between documents \mathbf{x}_i and \mathbf{x}_j as

$$\|\mathbf{x}_i - \mathbf{x}_j\| = \sqrt{\sum_{m=1}^J (x_{im} - x_{jm})^2}$$

Measuring the Distance Between Documents

Definition

The Euclidean distance between documents \mathbf{x}_i and \mathbf{x}_j as

$$\|\mathbf{x}_i - \mathbf{x}_j\| = \sqrt{\sum_{m=1}^J (x_{im} - x_{jm})^2}$$

Suppose $\mathbf{x}_i = (1, 4)$ and $\mathbf{x}_j = (2, 1)$. The distance between the documents is:

$$\begin{aligned}\|(1, 4) - (2, 1)\| &= \sqrt{(1 - 2)^2 + (4 - 1)^2} \\ &= \sqrt{10}\end{aligned}$$

Measuring the Distance Between Documents

Many distance metrics

Measuring the Distance Between Documents

Many distance metrics Consider the Minkowski family

Measuring the Distance Between Documents

Many distance metrics Consider the Minkowski family

Definition

The Minkowski Distance between documents \mathbf{X}_i and \mathbf{X}_j for value p is

$$d_p(\mathbf{X}_i, \mathbf{X}_j) = \left(\sum_{m=1}^J |x_{im} - x_{jm}|^p \right)^{1/p}$$

Members of the Minkowski Family

Members of the Minkowski Family

Manhattan metric

Members of the Minkowski Family

Manhattan metric

$$d_1(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^J |x_{im} - x_{jm}|$$

Members of the Minkowski Family

Manhattan metric

$$d_1(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^J |x_{im} - x_{jm}|$$
$$d_1((1, 4), (2, 1)) = |1| + |3| = 4$$

Members of the Minkowski Family

Manhattan metric

$$d_1(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^J |x_{im} - x_{jm}|$$
$$d_1((1, 4), (2, 1)) = |1| + |3| = 4$$

Minkowski (p) metric

Members of the Minkowski Family

Manhattan metric

$$d_1(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^J |x_{im} - x_{jm}|$$
$$d_1((1, 4), (2, 1)) = |1| + |3| = 4$$

Minkowski (p) metric

$$d_p(\mathbf{x}_i, \mathbf{x}_j) = \left(\sum_{m=1}^J |x_{im} - x_{jm}|^p \right)^{1/p}$$
$$d_p((1, 4), (2, 1)) = (|1 - 2|^p + |4 - 1|^p)^{1/p}$$

What Does p Do?

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences
If we let $p \rightarrow \infty$ Obtain maximum-metric (Chebyshev's Metric)

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences
If we let $p \rightarrow \infty$ Obtain maximum-metric (Chebyshev's Metric)

$$\lim_{p \rightarrow \infty} d_p(\mathbf{X}_i, \mathbf{X}_j) = \max_{m=1}^J |x_{im} - x_{jm}|$$

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences
If we let $p \rightarrow \infty$ Obtain maximum-metric (Chebyshev's Metric)

$$\lim_{p \rightarrow \infty} d_p(\mathbf{X}_i, \mathbf{X}_j) = \max_{m=1}^J |x_{im} - x_{jm}|$$

In words: distance between documents only the biggest difference

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences
If we let $p \rightarrow \infty$ Obtain maximum-metric (Chebyshev's Metric)

$$\lim_{p \rightarrow \infty} d_p(\mathbf{X}_i, \mathbf{X}_j) = \max_{m=1}^J |x_{im} - x_{jm}|$$

In words: distance between documents only the biggest difference
All other differences do not contribute to distance measure

What Does p Do?

Increasing $p \rightsquigarrow$ greater importance of coordinates with largest differences
If we let $p \rightarrow \infty$ Obtain maximum-metric (Chebyshev's Metric)

$$\lim_{p \rightarrow \infty} d_p(\mathbf{X}_i, \mathbf{X}_j) = \max_{m=1}^J |x_{im} - x_{jm}|$$

In words: distance between documents only the biggest difference
All other differences do not contribute to distance measure

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_2(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^2 + 4^2 + 3^2} = \sqrt{125} = 11.18$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_2(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^2 + 4^2 + 3^2} = \sqrt{125} = 11.18$$

$$d_4(\mathbf{X}_i, \mathbf{X}_j) = 10$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_2(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^2 + 4^2 + 3^2} = \sqrt{125} = 11.18$$

$$d_4(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_4(\mathbf{X}_i, \mathbf{X}_k) = \sqrt[4]{10^4 + 4^4 + 3^4} = (10337)^{1/4} = 10.08$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_2(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^2 + 4^2 + 3^2} = \sqrt{125} = 11.18$$

$$d_4(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_4(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^4 + 4^4 + 3^4} = (10337)^{1/4} = 10.08$$

$$d_\infty(\mathbf{X}_i, \mathbf{X}_j) = 10$$

Comparing the Metrics

Suppose $\mathbf{X}_i = (10, 4, 3)$, $\mathbf{X}_j = (0, 4, 3)$, and $\mathbf{X}_k = (0, 0, 0)$

Then:

$$d_1(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_1(\mathbf{X}_i, \mathbf{X}_k) = 10 + 4 + 3 = 17$$

$$d_2(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_2(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^2 + 4^2 + 3^2} = \sqrt{125} = 11.18$$

$$d_4(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_4(\mathbf{X}_i, \mathbf{X}_k) = \sqrt{10^4 + 4^4 + 3^4} = (10337)^{1/4} = 10.08$$

$$d_\infty(\mathbf{X}_i, \mathbf{X}_j) = 10$$

$$d_\infty(\mathbf{X}_i, \mathbf{X}_k) = 10$$

Are all differences equal?

Previous metrics treat all dimensions as **equal**

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Definition

Suppose that we have a covariance matrix Σ

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Definition

Suppose that we have a covariance matrix Σ . Then we can define the Mahalanobis Distance between documents \mathbf{X}_i and \mathbf{X}_j as ,

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Definition

Suppose that we have a covariance matrix Σ . Then we can define the Mahalanobis Distance between documents \mathbf{X}_i and \mathbf{X}_j as ,

$$d_{Mah}(\mathbf{X}_i, \mathbf{X}_j) = \sqrt{(\mathbf{X}_i - \mathbf{X}_j)' \Sigma^{-1} (\mathbf{X}_i - \mathbf{X}_j)}$$

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Definition

Suppose that we have a covariance matrix Σ . Then we can define the Mahalanobis Distance between documents \mathbf{X}_i and \mathbf{X}_j as ,

$$d_{Mah}(\mathbf{X}_i, \mathbf{X}_j) = \sqrt{(\mathbf{X}_i - \mathbf{X}_j)' \Sigma^{-1} (\mathbf{X}_i - \mathbf{X}_j)}$$

More generally: Σ could be symmetric and positive-definite

Are all differences equal?

Previous metrics treat all dimensions as **equal**

We may want to engage in some **scaling**/reweighting

Mahalanobis Distance

Definition

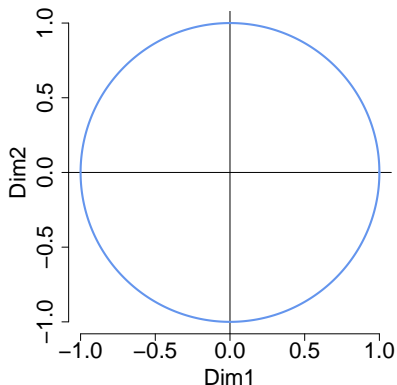
Suppose that we have a covariance matrix Σ . Then we can define the Mahalanobis Distance between documents \mathbf{X}_i and \mathbf{X}_j as ,

$$d_{Mah}(\mathbf{X}_i, \mathbf{X}_j) = \sqrt{(\mathbf{X}_i - \mathbf{X}_j)' \Sigma^{-1} (\mathbf{X}_i - \mathbf{X}_j)}$$

More generally: Σ could be symmetric and positive-definite

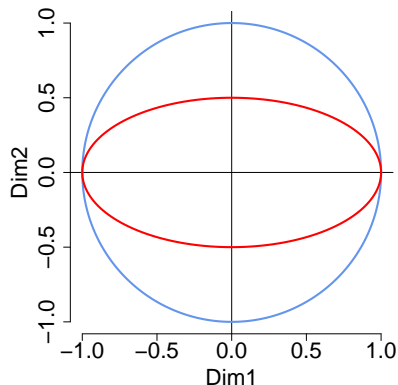
What does Σ do?

Some Intuition: The Unit Circle



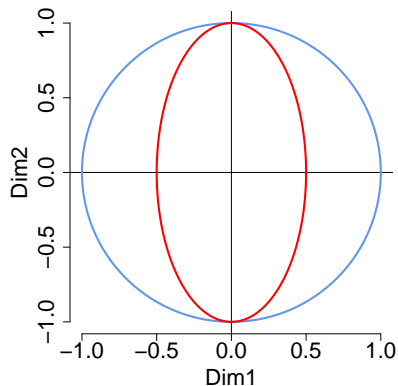
$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Some Intuition: The Unit Circle



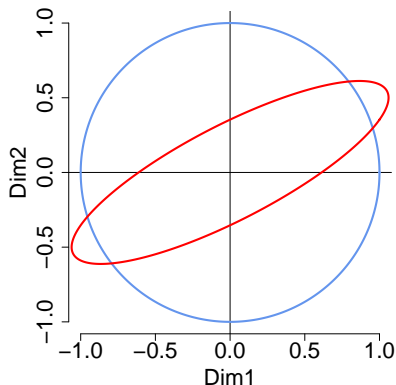
$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 0.5 \end{pmatrix}$$

Some Intuition: The Unit Circle



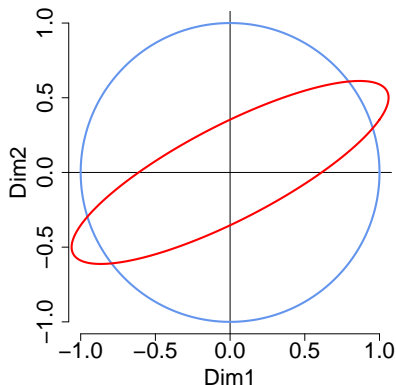
$$\Sigma = \begin{pmatrix} 0.5 & 0 \\ 0 & 1 \end{pmatrix}$$

Some Intuition: The Unit Circle



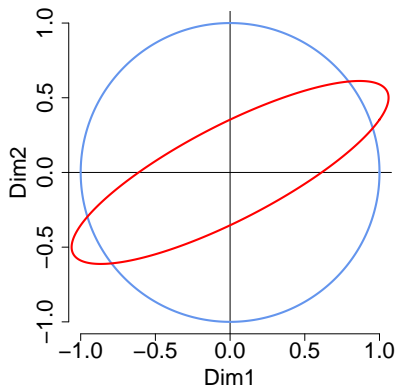
$$\Sigma = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 0.5 \end{pmatrix}$$

Some Intuition: The Unit Circle



$$\Sigma = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 0.5 \end{pmatrix}$$

Some Intuition: The Unit Circle



$$\Sigma = \begin{pmatrix} 1 & 0.3 \\ 0.3 & 0.5 \end{pmatrix}$$

Measuring Distance with Mahalanobis

Special Case 1: Identity Matrix

Measuring Distance with Mahalanobis

Special Case 1: Identity Matrix

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Measuring Distance with Mahalanobis

Special Case 1: Identity Matrix

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Then distance is **Euclidean**

Measuring Distance with Mahalanobis

Special Case 1: Identity Matrix

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Then distance is **Euclidean**

Special Case 2: Diagonal Matrix

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_J^2 \end{pmatrix}$$

Measuring Distance with Mahalanobis

Special Case 1: Identity Matrix

$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

Then distance is **Euclidean**

Special Case 2: Diagonal Matrix

$$\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_J^2 \end{pmatrix}$$

Then

$$d_{\text{Mah}}(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{m=1}^J \frac{(x_{im} - x_{jm})^2}{\sigma_m^2}}$$

Measuring Similarity

Measuring Similarity

What properties should similarity measure have?

Measuring Similarity

What properties should similarity measure have?

- Maximum: document with itself

Measuring Similarity

What properties should similarity measure have?

- Maximum: document with itself
- Minimum: documents have no words in common (**orthogonal**)

Measuring Similarity

What properties should similarity measure have?

- Maximum: document with itself
- Minimum: documents have no words in common (**orthogonal**)
- Increasing when **more** of same words used

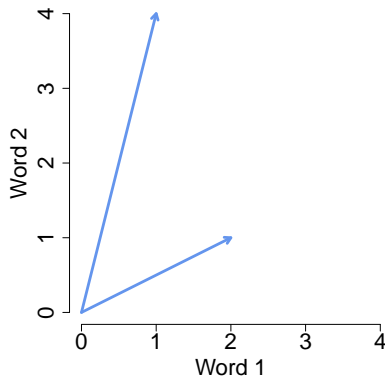
Measuring Similarity

What properties should similarity measure have?

- Maximum: document with itself
- Minimum: documents have no words in common (**orthogonal**)
- Increasing when **more** of same words used

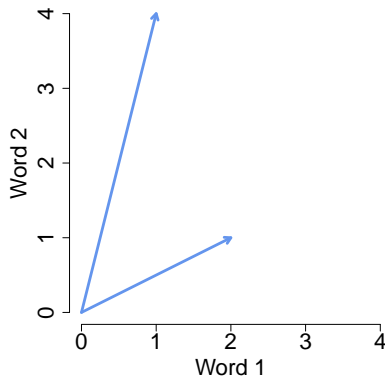
How should additional words be treated?

Measuring Similarity



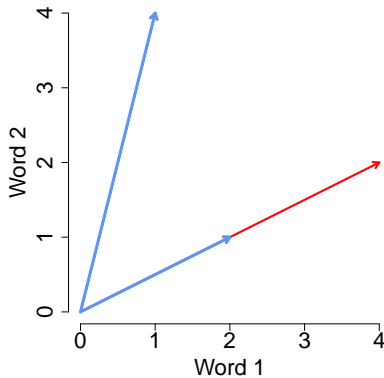
Measure 1: Inner product

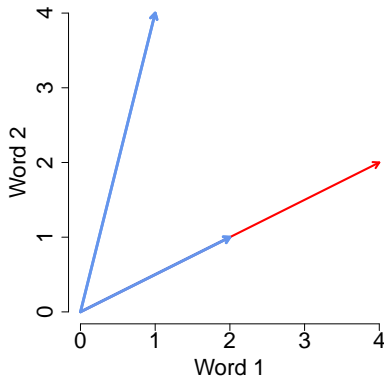
Measuring Similarity



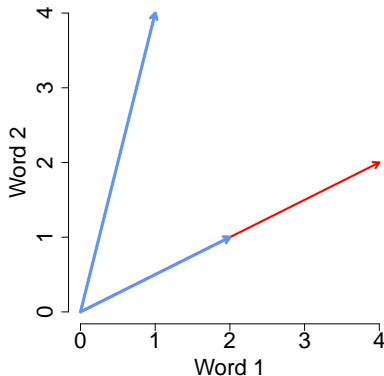
Measure 1: Inner product

$$(2, 1)' \cdot (1, 4) = 6$$



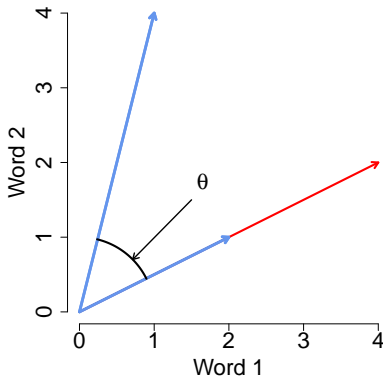


Problem(?): length dependent



Problem(?): length dependent

$$(4,2)'(1,4) = 12$$



Problem(?): length dependent

$$(4, 2)'(1, 4) = 12$$

$$a \cdot b = ||a|| \times ||b|| \times \cos \theta$$

Cosine Similarity

Cosine Similarity

$$\cos \theta = \left(\frac{a}{||a||} \right) \cdot \left(\frac{b}{||b||} \right)$$

Cosine Similarity

$$\cos \theta = \left(\frac{a}{\|a\|} \right) \cdot \left(\frac{b}{\|b\|} \right)$$

$$\frac{(4, 2)}{\|(4, 2)\|} = (0.89, 0.45)$$

Cosine Similarity

$$\cos \theta = \left(\frac{a}{\|a\|} \right) \cdot \left(\frac{b}{\|b\|} \right)$$

$$\frac{(4, 2)}{\|(4, 2)\|} = (0.89, 0.45)$$

$$\frac{(2, 1)}{\|(2, 1)\|} = (0.89, 0.45)$$

Cosine Similarity

$$\cos \theta = \left(\frac{a}{\|a\|} \right) \cdot \left(\frac{b}{\|b\|} \right)$$

$$\frac{(4, 2)}{\|(4, 2)\|} = (0.89, 0.45)$$

$$\frac{(2, 1)}{\|(2, 1)\|} = (0.89, 0.45)$$

$$\frac{(1, 4)}{\|(1, 4)\|} = (0.24, 0.97)$$

Cosine Similarity

$$\cos \theta = \left(\frac{a}{\|a\|} \right) \cdot \left(\frac{b}{\|b\|} \right)$$

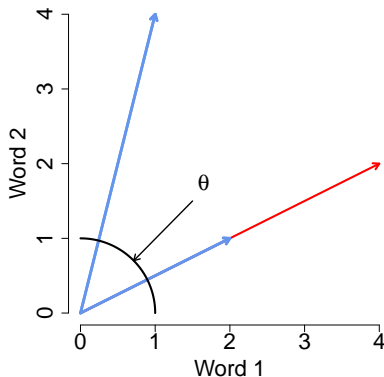
$$\frac{(4, 2)}{\|(4, 2)\|} = (0.89, 0.45)$$

$$\frac{(2, 1)}{\|(2, 1)\|} = (0.89, 0.45)$$

$$\frac{(1, 4)}{\|(1, 4)\|} = (0.24, 0.97)$$

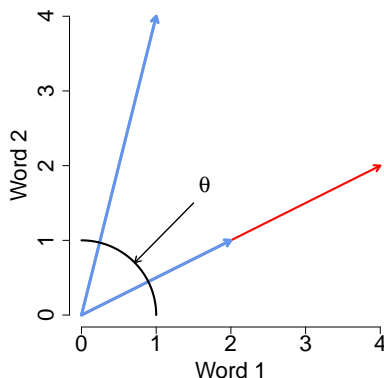
$$(0.89, 0.45)' (0.24, 0.97) = 0.65$$

Cosine Similarity



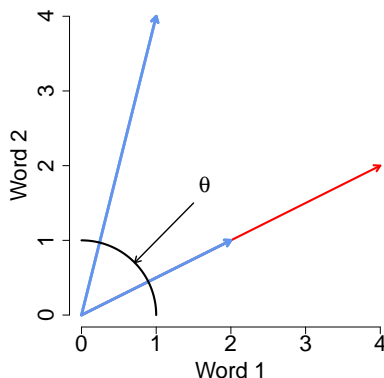
$\cos \theta$: removes document length from similarity measure

Cosine Similarity



$\cos \theta$: removes document length from similarity measure
Projects texts to unit length representation \rightsquigarrow onto sphere

Cosine Similarity



$\cos \theta$: removes document length from similarity measure
Projects texts to unit length representation \rightsquigarrow onto sphere

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\mathbf{x}_i^* \sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu})$$

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\begin{aligned}\mathbf{x}_i^* &\sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu}) \\ p(\mathbf{x}_i | \kappa, \boldsymbol{\mu}) &= c(\kappa) \exp(\kappa \mathbf{x}_i^* \boldsymbol{\mu})\end{aligned}$$

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\begin{aligned}\mathbf{x}_i^* &\sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu}) \\ p(\mathbf{x}_i | \kappa, \boldsymbol{\mu}) &= c(\kappa) \exp(\kappa \mathbf{x}_i^* \boldsymbol{\mu})\end{aligned}$$

Normal distribution, on a sphere

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\begin{aligned}\mathbf{x}_i^* &\sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu}) \\ p(\mathbf{x}_i | \kappa, \boldsymbol{\mu}) &= c(\kappa) \exp(\kappa \mathbf{x}_i^* \boldsymbol{\mu})\end{aligned}$$

Normal distribution, on a sphere

- Straightforward to Maximize

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\begin{aligned}\mathbf{x}_i^* &\sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu}) \\ p(\mathbf{x}_i | \kappa, \boldsymbol{\mu}) &= c(\kappa) \exp(\kappa \mathbf{x}_i^* \boldsymbol{\mu})\end{aligned}$$

Normal distribution, on a sphere

- Straightforward to Maximize
- Conjugate to itself

Von Mises-Fisher Distribution

Consider document \mathbf{x}_i .

$$\mathbf{x}_i^* = \frac{\mathbf{x}_i}{\|\mathbf{x}_i\|}$$

Then we might suppose:

$$\begin{aligned}\mathbf{x}_i^* &\sim \text{von Mises-Fisher}(\kappa, \boldsymbol{\mu}) \\ p(\mathbf{x}_i | \kappa, \boldsymbol{\mu}) &= c(\kappa) \exp(\kappa \mathbf{x}_i^* \boldsymbol{\mu})\end{aligned}$$

Normal distribution, on a sphere

- Straightforward to Maximize
- Conjugate to itself
- Useful for clustering, hierarchies of topics

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel of the *Gaussian* distribution

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel of the *Gaussian* distribution

σ^2 = determines sensitivity of the kernel

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel of the *Gaussian* distribution

σ^2 = determines sensitivity of the kernel

If $\mathbf{X}_i = \mathbf{X}_j$ then $k(\mathbf{X}_i, \mathbf{X}_j) = 1$

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel of the **Gaussian** distribution

σ^2 = determines sensitivity of the kernel

If $\mathbf{X}_i = \mathbf{X}_j$ then $k(\mathbf{X}_i, \mathbf{X}_j) = 1$

As \mathbf{X}_i and \mathbf{X}_j become more dissimilar, then $k(\mathbf{X}_i, \mathbf{X}_j) = 0$

Kernel Similarity

Definition

Suppose we have documents \mathbf{X}_i and \mathbf{X}_j . Define the *Gaussian* kernel as

$$k(\mathbf{X}_i, \mathbf{X}_j) = \exp\left(-\frac{\|\mathbf{X}_i - \mathbf{X}_j\|^2}{\sigma^2}\right)$$

Kernel of the *Gaussian* distribution

σ^2 = determines sensitivity of the kernel

If $\mathbf{X}_i = \mathbf{X}_j$ then $k(\mathbf{X}_i, \mathbf{X}_j) = 1$

As \mathbf{X}_i and \mathbf{X}_j become more dissimilar, then $k(\mathbf{X}_i, \mathbf{X}_j) = 0$

Result \rightsquigarrow often justify setting some kernel weights to zero

The Kernel Trick

Suppose all of our documents $\mathbf{x}_i \in \mathbb{R}^J$

The Kernel Trick

Suppose all of our documents $\mathbf{x}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

The Kernel Trick

Suppose all of our documents $\mathbf{x}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

The Kernel Trick

Suppose all of our documents $\mathbf{X}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

$$s(\phi(\mathbf{X}_i), \phi(\mathbf{X}_j)) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$$

The Kernel Trick

Suppose all of our documents $\mathbf{X}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

$$s(\phi(\mathbf{X}_i), \phi(\mathbf{X}_j)) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$$

- The only thing we care about, though is **inner product** of transformed variables

The Kernel Trick

Suppose all of our documents $\mathbf{X}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

$$s(\phi(\mathbf{X}_i), \phi(\mathbf{X}_j)) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$$

- The only thing we care about, though is **inner product** of transformed variables
- \rightsquigarrow So long as we can calculate inner product, we need not make explicit transformation

The Kernel Trick

Suppose all of our documents $\mathbf{X}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

$$s(\phi(\mathbf{X}_i), \phi(\mathbf{X}_j)) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$$

- The only thing we care about, though is **inner product** of transformed variables
- \rightsquigarrow So long as we can calculate inner product, we need not make explicit transformation
- \rightsquigarrow **Kernels** provide methods for capture wide array of transformations.

The Kernel Trick

Suppose all of our documents $\mathbf{X}_i \in \mathbb{R}^J$

There may be some mapping $\phi : \mathbb{R}^J \rightarrow \mathbb{R}^M$ where $M > J$ that improves our performance “lift” to higher dimension

We might want, then,

$$s(\phi(\mathbf{X}_i), \phi(\mathbf{X}_j)) = \langle \phi(\mathbf{X}_i), \phi(\mathbf{X}_j) \rangle$$

- The only thing we care about, though is **inner product** of transformed variables
- \rightsquigarrow So long as we can calculate inner product, we need not make explicit transformation
- \rightsquigarrow **Kernels** provide methods for capture wide array of transformations.
- **Kernel Trick** \rightsquigarrow calculate inner products on **untransformed** data (Gaussian Kernel), implicitly use wide array of ϕ 's.

Weighting Words

Are all words created equal?

Weighting Words

Are all words created equal?

- Treat all words equally

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words
 - Accentuate words that are likely to be informative

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words
 - Accentuate words that are likely to be informative
 - Make specific assumptions about characteristics of informative words

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words
 - Accentuate words that are likely to be informative
 - Make specific assumptions about characteristics of informative words

How to generate weights?

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words
 - Accentuate words that are likely to be informative
 - Make specific assumptions about characteristics of informative words

How to generate weights?

- Assumptions about separating words

Weighting Words

Are all words created equal?

- Treat all words equally
- Lots of noise
- Reweight words
 - Accentuate words that are likely to be informative
 - Make specific assumptions about characteristics of informative words

How to generate weights?

- Assumptions about separating words
- Use training set to identify separating words (Monroe, Ideology measurement)

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Ex. If all statements about OBL contain Bin Laden then this contributes nothing to similarity/dissimilarity measures

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Ex. If all statements about OBL contain Bin Laden then this contributes nothing to similarity/dissimilarity measures

Inverse document frequency:

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Ex. If all statements about OBL contain Bin Laden then this contributes nothing to similarity/dissimilarity measures

Inverse document frequency:

n_j = No. documents in which word j occurs

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Ex. If all statements about OBL contain Bin Laden then this contributes nothing to similarity/dissimilarity measures

Inverse document frequency:

n_j = No. documents in which word j occurs

$$\text{idf}_j = \log \frac{N}{n_j}$$

Weighting Words: TF-IDF Weighting

What properties do words need to separate concepts?

- Used frequently
- But not too frequently

Ex. If all statements about OBL contain Bin Laden then this contributes nothing to similarity/dissimilarity measures

Inverse document frequency:

n_j = No. documents in which word j occurs

$\text{idf}_j = \log \frac{N}{n_j}$

idf = $(\text{idf}_1, \text{idf}_2, \dots, \text{idf}_J)$

Weighting Words: TF-IDF Weighting

Why log ?

Weighting Words: TF-IDF Weighting

Why log ?

- Maximum at $n_j = 1$

Weighting Words: TF-IDF Weighting

Why log ?

- Maximum at $n_j = 1$
- Decreases at rate $\frac{1}{n_j} \Rightarrow$ diminishing “penalty” for more common use

Weighting Words: TF-IDF Weighting

Why log ?

- Maximum at $n_j = 1$
- Decreases at rate $\frac{1}{n_j} \Rightarrow$ diminishing “penalty” for more common use
- Other functional forms are fine, embed assumptions about penalization of common use

Weighting Words: TF-IDF

Weighting Words: TF-IDF

$$\mathbf{X}_{i,\text{idf}} \equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J)$$

Weighting Words: TF-IDF

$$\mathbf{X}_{i,\text{idf}} \equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J)$$
$$\mathbf{X}_{j,\text{idf}} \equiv \mathbf{X}_j \times \mathbf{idf} = (X_{j1} \times \text{idf}_1, X_{j2} \times \text{idf}_2, \dots, X_{jJ} \times \text{idf}_J)$$

Weighting Words: TF-IDF

$$\mathbf{X}_{i,\text{idf}} \equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J)$$
$$\mathbf{X}_{j,\text{idf}} \equiv \mathbf{X}_j \times \mathbf{idf} = (X_{j1} \times \text{idf}_1, X_{j2} \times \text{idf}_2, \dots, X_{jJ} \times \text{idf}_J)$$

How Does This Matter For Measuring Similarity/Dissimilarity?

Weighting Words: TF-IDF

$$\mathbf{X}_{i,\text{idf}} \equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J)$$
$$\mathbf{X}_{j,\text{idf}} \equiv \mathbf{X}_j \times \mathbf{idf} = (X_{j1} \times \text{idf}_1, X_{j2} \times \text{idf}_2, \dots, X_{jJ} \times \text{idf}_J)$$

How Does This Matter For Measuring Similarity/Dissimilarity?

Inner Product

Weighting Words: TF-IDF

$$\mathbf{X}_{i,\text{idf}} \equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J)$$
$$\mathbf{X}_{j,\text{idf}} \equiv \mathbf{X}_j \times \mathbf{idf} = (X_{j1} \times \text{idf}_1, X_{j2} \times \text{idf}_2, \dots, X_{jJ} \times \text{idf}_J)$$

How Does This Matter For Measuring Similarity/Dissimilarity?

Inner Product

$$\mathbf{X}_{i,\text{idf}} \cdot \mathbf{X}_{j,\text{idf}} = (\mathbf{X}_i \times \mathbf{idf})' (\mathbf{X}_j \times \mathbf{idf})$$

Weighting Words: TF-IDF

$$\begin{aligned}\mathbf{X}_{i,\text{idf}} &\equiv \underbrace{\mathbf{X}_i}_{\text{tf}} \times \mathbf{idf} = (X_{i1} \times \text{idf}_1, X_{i2} \times \text{idf}_2, \dots, X_{iJ} \times \text{idf}_J) \\ \mathbf{X}_{j,\text{idf}} &\equiv \mathbf{X}_j \times \mathbf{idf} = (X_{j1} \times \text{idf}_1, X_{j2} \times \text{idf}_2, \dots, X_{jJ} \times \text{idf}_J)\end{aligned}$$

How Does This Matter For Measuring Similarity/Dissimilarity?

Inner Product

$$\begin{aligned}\mathbf{X}_{i,\text{idf}} \cdot \mathbf{X}_{j,\text{idf}} &= (\mathbf{X}_i \times \mathbf{idf})' (\mathbf{X}_j \times \mathbf{idf}) \\ &= (\text{idf}_1^2 \times X_{i1} \times X_{j1}) + (\text{idf}_2^2 \times X_{i2} \times X_{j2}) + \\ &\quad \dots + (\text{idf}_J^2 \times X_{iJ} \times X_{jJ})\end{aligned}$$

Weighting Words: Inner Product

Define:

Weighting Words: Inner Product

Define:

$$\Sigma = \begin{pmatrix} \text{idf}_1^2 & 0 & 0 & \dots & 0 \\ 0 & \text{idf}_2^2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \text{idf}_J^2 \end{pmatrix}$$

Weighting Words: Inner Product

Define:

$$\Sigma = \begin{pmatrix} \text{idf}_1^2 & 0 & 0 & \dots & 0 \\ 0 & \text{idf}_2^2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \text{idf}_j^2 \end{pmatrix}$$

If we use tf-idf for our documents, then

Weighting Words: Inner Product

Define:

$$\Sigma = \begin{pmatrix} \text{idf}_1^2 & 0 & 0 & \dots & 0 \\ 0 & \text{idf}_2^2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \text{idf}_J^2 \end{pmatrix}$$

If we use tf-idf for our documents, then

$$\begin{aligned} d_2(\mathbf{x}_i, \mathbf{x}_j) &= \sqrt{\sum_{m=1}^J (x_{im,\text{idf}} - x_{jm,\text{idf}})^2} \\ &= \sqrt{(\mathbf{x}_i - \mathbf{x}_j)' \Sigma (\mathbf{x}_i - \mathbf{x}_j)} \end{aligned}$$

Final Product

Applying some measure of distance, similarity (if symmetric) yields:

$$\mathbf{D} = \begin{pmatrix} 0 & d(1,2) & d(1,3) & \dots & d(1,N) \\ d(2,1) & 0 & d(2,3) & \dots & d(2,N) \\ d(3,1) & d(3,2) & 0 & \dots & d(3,N) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ d(N,1) & d(N,2) & d(N,3) & \dots & 0 \end{pmatrix}$$

Lower Triangle contains unique information $N(N-1)/2$

R Code!!

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

Spiraling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans

Why?

- American political development

Spiraling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- **Political Science question**: how did Native Americans lose land so quickly?

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- **Political Science question**: how did Native Americans lose land so quickly?

Paper does **a lot**. We're going to focus on

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- **Political Science question**: how did Native Americans lose land so quickly?

Paper does **a lot**. We're going to focus on

- Today: Text representation and similarity calculation

Spirling and Indian Treaties

Spirling (2013): model Treaties between US and Native Americans
Why?

- American political development
- IR Theories of Treaties and Treaty Violations
- Comparative studies of indigenous/colonialist interaction
- **Political Science question**: how did Native Americans lose land so quickly?

Paper does **a lot**. We're going to focus on

- Today: Text representation and similarity calculation
- Tuesday: Projecting to low dimensional space

Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between Us

Analyzes K-substrings

Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between Us

Analyzes K-substrings

Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between Us

Analyzes K-substrings

Spiraling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spiraling uses complicated representation of texts to preserve word order ~→
broad application

Peace **B**etween Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace **B**etween Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace **Between** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace **Between** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Bet**ween** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Bet**ween** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between **een** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between **en** Us

Analyzes K-substrings

Spirling and Indian Treaties

How do we preserve word order and semantic language?

After stemming, stopping, bag of wording:

- Peace Between Us
- No Peace Between Us

are identical.

Spirling uses complicated representation of texts to preserve word order ~→
broad application

Peace Between **Us**

Analyzes K-substrings

Kernel Trick

Kernel Trick

- **Kernel Methods:** Represent texts, measure similarity

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Uses Kernel methods to measure **similarity**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Uses Kernel methods to measure **similarity**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Uses Kernel methods to measure **similarity**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Uses Kernel methods to measure **similarity**

Kernel Trick

- **Kernel Methods**: Represent texts, measure similarity **simultaneously**
- Compare only **substrings** in both documents (without explicitly quantifying entire documents)
- Problem solved:
 - **Arthur** gives all his money to **Justin**
 - **Justin** gives all his money to **Arthur**
 - Discard word order: same sentence Kernel : different sentences.

Uses Kernel methods to measure **similarity**

Similarity and Dissimilarity of Many Things

Throughout the course we'll measure **similarity** between documents
We'll also (implicitly) study **similarity of probability distributions**
Develop a measure of distribution dissimilarity

Similarity of Probability Distributions

Definition

Suppose P is a continuous random variable with density $p : \mathbb{R} \rightarrow \mathbb{R}$ and Q is a continuous random variable with density $q : \mathbb{R} \rightarrow \mathbb{R}$.

We can define the KL-Divergence between P and Q as

$$KL(P||Q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

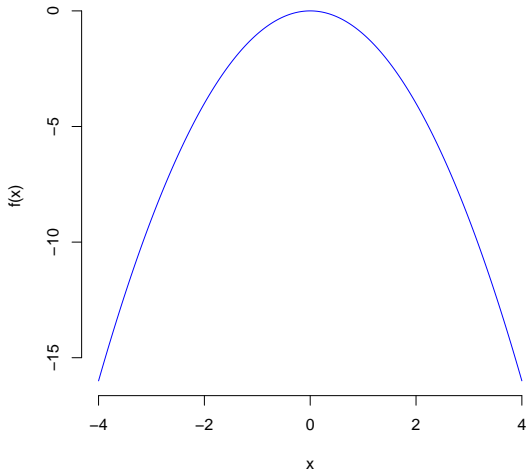
Assessing Similarity of Other Things

KL-divergence measures **dissimilarity** between two distributions.

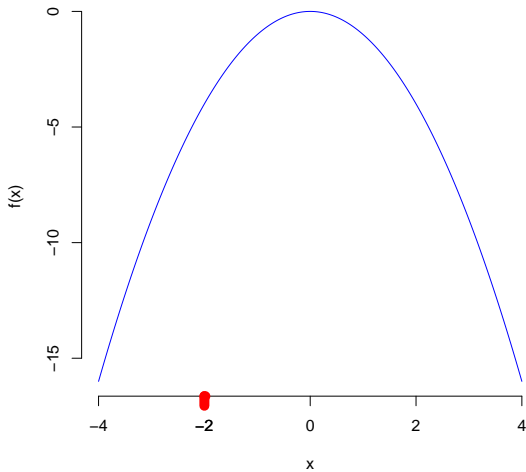
Consider a function. $f(x) = -x^2$.

Consider a function. $f(x) = -x^2$.
Maps numbers to other numbers.

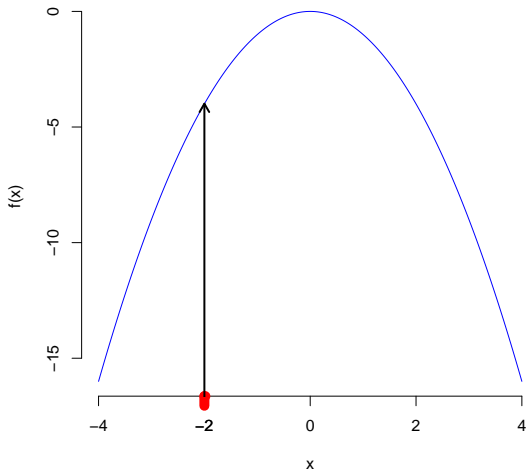
Consider a function. $f(x) = -x^2$.
Maps numbers to other numbers.



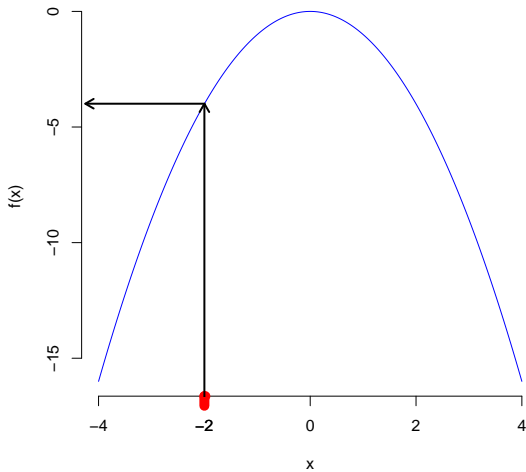
Take some input (-2 here)



Then obtain the value of $f(-2)$



Then obtain the value of $f(-2) = -4$



$\text{KL}(q||p)$ is a **functional**.

$KL(q||p)$ is a **functional**. A functional takes **functions** as inputs, returns a real number.

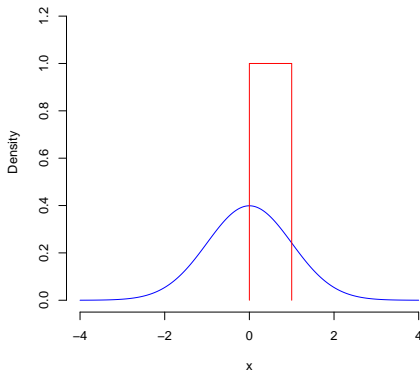
$\text{KL}(q||p)$ is a **functional**. A functional takes **functions** as inputs, returns a real number.

$\text{KL}(q||p)$ maps from sets of distributions $q \in \mathcal{Q}$ and $p \in \mathcal{P}$ to positive real numbers.

$KL(q||p)$ is a **functional**. A functional takes **functions** as inputs, returns a real number.

$KL(q||p)$ maps from sets of distributions $q \in \mathcal{Q}$ and $p \in \mathcal{P}$ to positive real numbers.

For example, we could set $q = \text{Uniform}(0,1)$ and $p = \text{Normal}(0, 1)$

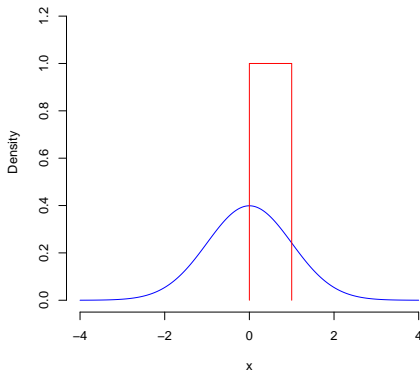


$KL(q||p)$ is a **functional**. A functional takes **functions** as inputs, returns a real number.

$KL(q||p)$ maps from sets of distributions $q \in \mathcal{Q}$ and $p \in \mathcal{P}$ to positive real numbers.

For example, we could set $q = \text{Uniform}(0,1)$ and $p = \text{Normal}(0, 1)$

$$KL(\text{Uniform}(0,1)||\text{Normal}(0,1)) = 1.09$$



If q and p are the **same** distribution then $\text{KL}(q||p) = 0$.

If q and p are the **same** distribution then $KL(q||p) = 0$.

Variational Approximation (topic models!): **approximate** one distribution p , with another, simpler distribution q .

If q and p are the **same** distribution then $KL(q||p) = 0$.

Variational Approximation (topic models!): **approximate** one distribution p , with another, simpler distribution q .

Then make this approximation the **best** possible—minimize the KL-divergence.

A simple example.

A simple example.

Approximate a $\text{Normal}(0,1)$ with symmetric Uniform distribution,
 $\text{Uniform}(-b, b)$.

A simple example.

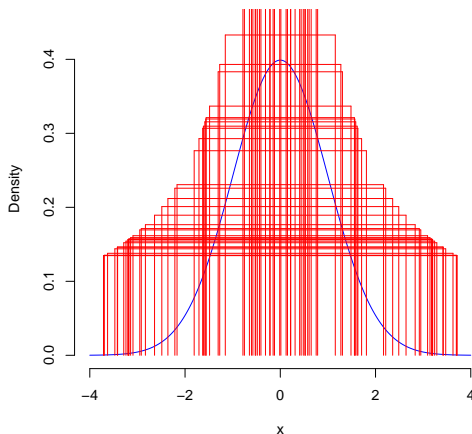
Approximate a $\text{Normal}(0,1)$ with symmetric Uniform distribution, $\text{Uniform}(-b, b)$.

Choose b to min. $\text{KL}(\text{Uniform}(-b, b) \parallel \text{Normal}(0,1))$

A simple example.

Approximate a $\text{Normal}(0,1)$ with symmetric $\text{Uniform}(-b, b)$.

Choose b to min. $\text{KL}(\text{Uniform}(-b, b) || \text{Normal}(0,1))$



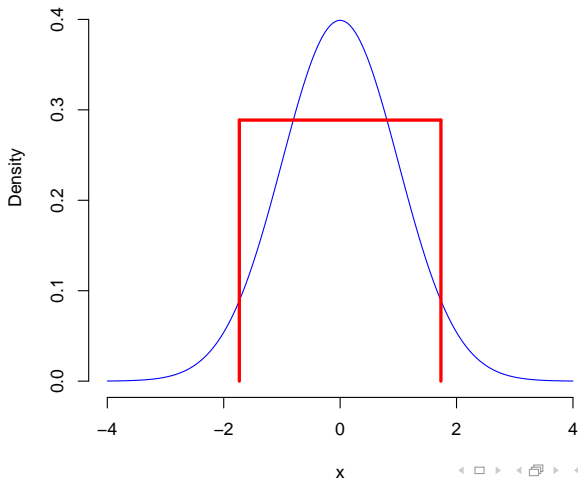
Answer:

Answer:

$$b = \sqrt{3}$$

Answer:

$$b = \sqrt{3}$$



- 1) Documents in vector space \rightsquigarrow geometry of texts
- 2) Many methods to measure similarity and dissimilarity