

Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selection

Wrappers

Embedded

Filtering method

PMI MI

Common lesi

Practical

Summar

## **Lecture 18: Feature Selection**

## COMP90049 Knowledge Technologies

Sarah Erfani and Karin Verspoor, and Jeremy Nicholson, CIS

Semester 2, 2018





#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Footure Coloctic

Wrappers Embedded

Filtering method

Filtering method

MI ~2

Common Issues

Practical consideration

Summary



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selection

Embedded Filters

Filtering method

PMI

 $\chi^2$ 

Practical consideration

Summary

## Data set:

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
:	:	:	:	:



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

### Features in Machine Learning

Feature Select

Wrappers Embedded

Filters

Filtering method

MI

X

Oommon 133uc

Practical consideration

Summary

## Instances:

Outlook	Temperature	Humidity	Windy	Play
surny	Sot T		TALE	<del>100</del> 1
suiny	Sot T		TRE	$\frac{1}{2}$
overcast	hot	high	FALSE	yes 2
rainy	mild	high	FALSE	yes
rainy	cool	n or m a l	FALSE	yes
rainy	cool	normal	TRUE	no
:	:	:	:	:



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

## Features in Machine Learning

Wrappers
Embedded

Filters

Filtering method

MI

*x* -

Common issue

Practical consideration

Summary

## Attributes:

Outlook	Temperature	Humidity	Windy	Play
sunny	Hot	high	FALSE	no
suAhy	i <b>Z</b>	high	TRUE	no
overcast	hor	high	FALSE	yes
ra	m <mark>H</mark> d	high	FALSE	yes
ra <del>in</del> y	c <del>oq</del> 1	n or m a l	FALSE	yes
ra <del>in</del> y	c <del>log</del> l	normal	TRUE	no
I	Ţ	:	:	:
<del></del>	12			



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Select

Wrappers Embedded

Embedded Filters

Filtering method

PMI

MI

Common Issues

Practical

Summar

- Where do instances come from?
  - Examples from real world data



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Select

Wrappers Embedded

Embedded Filters

Filtering method

PMI

Commo

Common issues

Practical consideration

Summary

- Where do instances come from?
  - Examples from real world data
- Where do attributes come from?
  - **???**



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

**Feature Selection** 

Wrappers Embedded

Filters

Filtering method

PMI MI

Common Iss

Practical consideration

Summary

- Where do instances come from?
  - Examples from real world data
- Where do attributes come from?
  - (Hopefully) meaningful features of the problem
  - Anything that might capture regularity in the data



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selecti

Wrappers Embedded

Filters

Filtering method

PMI

MI

Common Issue

Practical consideration

Summary

Outlook	Temperature	Humidity	Windy	Play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
:	<b>:</b>	÷	:	:

- Windy seems like a good predictor of Play
- Humidity seems like a less good predictor of Play



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Embedded Filters

Filtering method

PMI MI

Common Is

Practical

Summary

- Where do instances come from?
  - Examples from real world data
- Where do attributes come from?
  - (Hopefully) meaningful features of the problem
  - Anything that might capture regularity in the data
- Where do models come from?
  - **???**



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Select

Wrappers Embedded

Filters

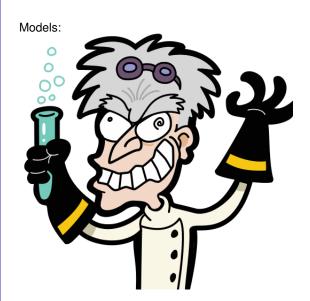
Filtering metho

MI

O-------

Practical

Summary





### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection Wrappers

Embedded Filters

Filtering method PMI

 $\frac{\mathsf{MI}}{\chi^2}$ 

Common issues

Practical consideration

Summary

- Where do instances come from?
  - Examples from real world data
- Where do attributes come from?
  - (Hopefully) meaningful features of the problem
  - Anything that might capture regularity in the data
- Where do models come from?
  - Need to choose a model suitable for our data set



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Wrappers Embedded

ilters

Filtering method

PMI MI

 $\chi^2$ 

Common Issues

Practical consideration

.....

- Pick a feature representation
- 2
- 3
- 4
- 5
- 6



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Wrappers Embedded

Embedded Filters

Filtering method

PMI

MI

Common Iss

Practical consideration

.....

- Pick a feature representation
- Compile data
- 3
- 4
- 5
- 6



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Embedded

ilters

Filtering method

PMI

 $\chi^2$ 

Common Issues

Practical consideration

.....

- Pick a feature representation
- Compile data
- 3 Pick a (suitable) algorithm for building a model
- 4
- 5
- 6



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Embedded

ilters

Filtering methods

PMI

MI

Common Issu

Practical consideration

Lummaru

- Pick a feature representation
- Compile data
- Pick a (suitable) algorithm for building a model
- Train the model
- 5
- 6



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Embedded

Filters

Filtering method

PMI

 $\chi^{^2}$ 

Common Issues

Practical consideration

Summar

- Pick a feature representation
- Compile data
- 3 Pick a (suitable) algorithm for building a model
- Train the model
- Classify development data, evaluate results
- 6



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Wrappers Embedded

Filters

Filtering method

PMI

MI

Common Issu

Practical consideration

Summar

- Pick a feature representation
- Compile data
- 3 Pick a (suitable) algorithm for building a model
- Train the model
- Classify development data, evaluate results
- Probably: go to (1)



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection

Wrappers

Filters

Filtering method

PMI

MI

Common Issues

Practical consideration

Summar

- [0.] Get hired!
- Pick a feature representation
- Compile data
- Pick a (suitable) algorithm for building a model
- 5 Train the model
- Classify development data, evaluate results
- Probably: go to (1)



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selecti

Wrappers Embedded

Filtering method

PMI MI

Common Issues

Practical consideration

Summary

- Choose an algorithm suitable for classifying the data according to the attributes
- Choose attributes suitable for classifying the data according to the model



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selecti

Wrappers Embedded

Filtering method

PMI MI

Common Issue

Practical consideration

Summar

- Choose an algorithm suitable for classifying the data according to the attributes
- Choose attributes suitable for classifying the data according to the model
  - Inspection
    - Intuition



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Wrappers Embedded

Filtering method

PMI MI

Common Issi

Practical consideration

Summary

- Choose an algorithm suitable for classifying the data according to the attributes
- Choose attributes suitable for classifying the data according to the model
  - Inspection
  - Intuition
  - Neither possible in practice



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

#### Features in Machine Learning

Feature Selection
Wrappers

Embedded Filters

Filtering method

 $\chi^2$ 

Common issues

Practical consideration

Summary

- Choose an algorithm suitable for classifying the data according to the attributes
- Choose attributes suitable for classifying the data according to the model
  - Inspection
  - Intuition
  - Neither possible in practice
  - Throw everything we can think of at the problem and let the algorithm decide!



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in

# Feature Selection

Wrappers Embedded

Filtering methods

PMI MI

Common Issues

Practical consideration

Summary



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Feature Selection

Wrappers Embedded

Filters

Filtering methods

PMI MI

Common Issues

Practical

Summarı

### Better models!

■ Better performance according to some evaluation metric



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### **Feature Selection**

Wrappers Embedded Filters

Filtering method

PMI

Common Issue

Practical consideration

Summary

### Better models!

Better performance according to some evaluation metric

## Side-goal:

■ Tell us interesting things about the problem



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

#### **Feature Selection**

Wrappers Embedded Filters

Filtering methods

PMI MI

Common Issues

Practical consideration

Summa

### Better models!

Better performance according to some evaluation metric

## Side-goal:

■ Tell us interesting things about the problem

## Side-goal:

Fewer features → smaller models → faster answer



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

#### **Feature Selection**

Embedded Filters

Filtering method

PMI MI

Common Issu

Practical consideration

Summar

### Better models!

Better performance according to some evaluation metric

## Side-goal:

■ Tell us interesting things about the problem

## Side-goal:

- Fewer features → smaller models → faster answer
  - More accurate answer >> faster answer



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

### Wrappers

Embedded

Filtering method

PMI

 $\chi^{^2}$ 

Common Issues

Practical consideration

Summary

## "Wrapper" methods:

 Choose subset of attributes that give best performance on the development data



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

Feature Selection

### Wrappers

Embedded

Filtering method

PMI

Common Iss

Practical

Summary

- Choose subset of attributes that give best performance on the development data
- For example: for the Weather data set:
  - Train model on {Outlook}



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Feature Selec

## Wrappers

Filters

Filtering method

MI

Common Issue

Practical consideration

Summary

- Choose subset of attributes that give best performance on the development data
- For example: for the Weather data set:
  - Train model on {Outlook}
  - Train model on {Temperature}
    - ---



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning
Feature Selection

Wrappers Embedded

Filtering method

PMI MI

Common issue

Practical consideration

Summary

- Choose subset of attributes that give best performance on the development data
- For example: for the Weather data set:
  - Train model on {Outlook}
  - Train model on {Temperature}
  - Train model on {Outlook, Temperature}

  - Train model on {Outlook, Temperature, Humidity}
  - ...
  - Train model on {Outlook, Temperature, Humidity, Windy}



## Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning
Feature Selection

Wrappers Embedded

Filters

Filtering method

Common Issue

Practical consideration

Summary

- Choose subset of attributes that give best performance on the development data
- For example: for the Weather data set:
  - Evaluate model on {Outlook}
  - Evaluate model on {Temperature}
  - Evaluate model on {Outlook, Temperature}
  - ...Evaluate model on {Outlook, Temperature, Humidity}
  - ...
  - Evaluate model on {Outlook, Temperature, Humidity, Windy}



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selection

Wrappers Embedded

Filters

Filtering moth

PMI

X Common les

Practical

onsideration

- Choose subset of attributes that give best performance on the development data
- For example: for the Weather data set:
  - Evaluate model on {Outlook}
  - Evaluate model on {Temperature}
    - ...
  - Evaluate model on {Outlook, Temperature}
  - ...
  - Evaluate model on {Outlook, Temperature, Humidity}
  - Evaluate model on {Outlook, Temperature, Humidity, Windy}
- Best performance on data set → best feature set



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Feature Selection

Wrappers Embedded

Filters

Filtering method

PMI MI

Common Issues

Practical consideration

Summary

- Choose subset of attributes that give best performance on the development data
- Advantages:
  - Feature set with optimal performance on development data



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

### Wrappers

rinters

PMI

MI

Common Issues

Practical consideration

Summary

- Choose subset of attributes that give best performance on the development data
- Advantages:
  - Feature set with optimal performance on development data
- Disadvantages:
  - Takes a long time



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learni

### Feature Selecti

Wrappers Embedded

Filtering method

PMI MI

Common Issue

Practical consideration

Summar

Assume we have a fast method (e.g. Naive Bayes) over a data set of non-trivial size ( $\sim$ 10K instances):

Assume: train-evaluate cycle takes 10 sec to complete



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

**Feature Selection** 

Wrappers Embedded

Filtering method

PMI MI

Common Issues

Practical consideration

Summary

Assume we have a fast method (e.g. Naive Bayes) over a data set of non-trivial size ( $\sim$ 10K instances):

Assume: train-evaluate cycle takes 10 sec to complete

How many cycles? For *m* features:

■ 
$$2^m$$
 subsets =  $\frac{2^m}{6}$  minutes



### Lecture 18: Feature Selection

COMP90049 Knowledge **Technologies** 

Machine Learning

Wrappers

Assume we have a fast method (e.g. Naive Bayes) over a data set of non-trivial size ( $\sim$ 10K instances):

Assume: train-evaluate cycle takes 10 sec to complete

How many cycles? For *m* features:

- $\blacksquare$  2<sup>m</sup> subsets =  $\frac{2^m}{6}$  minutes
- $= m = 10 \rightarrow 3 \text{ hours}$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Lear

Feature Selection
Wrappers

Embedded

Filtering moth

PMI

 $\chi^2$ 

Common Issues

Practical consideration

Summar

Assume we have a fast method (e.g. Naive Bayes) over a data set of non-trivial size ( $\sim$ 10K instances):

Assume: train-evaluate cycle takes 10 sec to complete

How many cycles? For *m* features:

- $2^m$  subsets =  $\frac{2^m}{6}$  minutes
- $= m = 10 \rightarrow 3 \text{ hours}$
- $= m = 60 \rightarrow \text{heat death of universe}$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learn

Feature Selection

Embedded Filters

Filtering method

PMI MI

Common Iss

Practical consideration

Summar

Assume we have a fast method (e.g. Naive Bayes) over a data set of non-trivial size ( $\sim$ 10K instances):

Assume: train-evaluate cycle takes 10 sec to complete

How many cycles? For *m* features:

- $2^m$  subsets =  $\frac{2^m}{6}$  minutes
- $m = 10 \rightarrow 3$  hours
- $= m = 60 \rightarrow \text{heat death of universe}$

Only practical for very small data sets.



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### . . . . . .

#### Wrappers

Embedded

Filtering method

### PMI

 $\chi^2$ 

Common Issues

Practical consideration

Summary

- Train and evaluate model on each single attribute
- Choose best attribute



### Lecture 18: Feature Selection

COMP90049 Knowledge **Technologies** 

Machine Learning

Feature Selection Wrappers

- Train and evaluate model on each single attribute
- Choose best attribute
- Until convergence:
  - Train and evaluate model on best attribute(s), plus each remaining single attribute
  - Choose best attribute out of the remaining set
- Iterate until performance (e.g. accuracy) stops increasing



#### Lecture 18: **Feature Selection**

COMP90049 Knowledge **Technologies** 

Wrappers

- Bad news:
  - Takes ½m² cycles, for m attributes
     In theory, 386 attributes → days



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Feature Selection

Wrappers

Embedded

Filtering method

DMI

MI

Common Issi

Practical consideration

Summary

- Bad Good news:
  - Takes  $\frac{1}{2}m^2$  cycles, for m attributes
  - In practice, converges much more quickly than this



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Feature Selection

### Wrappers

Eiltore

Filtering method

DMI

PMI MI

Common iss

Practical consideration

Summary

- Bad Good Bad news:
  - Takes  $\frac{1}{2}m^2$  cycles, for m attributes
  - In practice, converges much more quickly than this
  - Convergences to a sub-optimal (and often very bad) solution



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

### Wrappers

Embedded

Filtering method

PMI method

MI

Common Issue

Practical consideration

Summary

- Start with all attributes
- Remove one attribute, train and evaluate model



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection
Wrappers

Embedded

Filters

Filtering method

Common Issu

Practical

Summar

- Start with all attributes
- Remove one attribute, train and evaluate model
- Until divergence:
  - From remaining attributes, remove each attribute, train and evaluate model
  - Remove attribute that causes least performance degredation
- $\blacksquare$  Termination condition usually: performance (e.g. accuracy) starts to degrade by more than  $\epsilon$



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Eastura Calastia

#### Wrappers

Filters

Filtering method

### PMI

MI 2

Common Issues

Practical

Summary

- Good news:
  - Mostly removes irrelevant attributes (at the start)



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

**Feature Selection** 

#### Wrappers

Filters

Filtering method

Filtering method

PMI

~

Common Issues

Practical consideration

Summary

- Good news:
  - Mostly removes irrelevant attributes (at the start)
- Bad news:
  - Assumes independence of attributes (both approaches; worse than Naive Bayes!)
  - Actually does take  $O(m^2)$  time; cycles are slower with more attributes



### In-built feature selection

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Heatures in Machine Learnin

Feature Selection

Embedded

Filtering method

PMI

Common Issue

Practical

Summary

### "Embedded" methods:

Some models actually perform feature selection as part of the algorithm!



### In-built feature selection

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Wrappers

Embedded

Filtering method

PMI

Common Is

Practical consideration

Summary

#### "Embedded" methods:

- Some models actually perform feature selection as part of the algorithm!
  - Most notably, linear classifiers
  - To some degree: Decision Trees
- (More on these models in later lectures.)



### In-built feature selection

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection
Wrappers

Embedded

Filtering method

PMI

 $\chi^2$ 

Common Issues

Practical consideration

Summary

#### "Embedded" methods:

- Some models actually perform feature selection as part of the algorithm!
  - Most notably, linear classifiers
  - To some degree: Decision Trees
- Often benefit from other feature selection approaches anyway



# Feature filtering

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnir

#### E. ........

Wrappers

#### Embedde Filters

Filtering method

PMI

MI

Common Issues

Practical consideration

Summar

Intuition: possible to evaluate "goodness" of each feature, separate from other features



# Feature filtering

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Select

Embedde

Filters

Filtering method PMI

X

Practical

Summary

Intuition: possible to evaluate "goodness" of each feature, separate from other features

- Consider each feature separately: linear time in number of attributes
- Typically most popular strategy
- Possible (but difficult) to control for inter-dependence of features



#### Lecture 18: **Feature Selection**

COMP90049 Technologies

Filtering methods

What makes a feature set good?



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### \_\_\_\_\_

Embedded

Filters

#### Filtering methods

MI 2

Common Issue

Practical

Summar

What makes a feature set good?

■ Better models!



#### Lecture 18: **Feature Selection**

COMP90049

Filtering methods

What makes a feature set single feature good?

Better models!



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Foature Salactic

Embedded Filters

Filtering methods

#### PMI

MI

Common Issues

Practical

Summary

What makes a feature set single feature good?

- Better models!
- Well correlated with class



# Toy example

Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnin

Footure Coloctics

Wrappers

Filters

Filtering methods

МІ

Commo

Practical

consideratior

Summarv

$a_1$	$a_2$	С
Y	Υ	Υ
Υ	Ν	Υ
Ν	Υ	Ν
Ν	N	Ν

Which of  $a_1$ ,  $a_2$  is good?



# Toy example

Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

E. ........

Wrappers

Filters

Filtering methods

MI

х -

Common issue

Practical consideration

Summary

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	N
Ν	Ν	N

 $a_1$  is probably good.



# Toy example

Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnin

Factions Calcution

Wrappers Embedde

Filters

Filtering methods

МІ

Commo

0011111011110000

Practical consideration

Summary

$a_1$	$a_2$	C
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	N
Ν	Ν	N

 $a_2$  is probably not good.



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Factions Calcution

Wrappers Embedded

Embedded Filters

Filtering metho

PMI

Common

Practical consideration

Summary

### Recall independence:

$$P(A,C)=P(A)P(C)$$



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnii

#### Eggtura Calagtic

Embedded

Filtering method

# PMI

MI

Common Issues

Practical consideration

Summary

Recall independence:

$$P(A, C) = P(A)P(C)$$

This formula holds if attribute is independent from class.



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

Feature Selection

Wrappers Embedded

Embedded Filters

Filtering method

PMI

χ

Common issue:

Practical consideration

Summary

## Recall independence:

$$\frac{P(A,C)}{P(A)P(C)}=1$$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning
Feature Selection

Wrappers Embedded

Embedded Filters

Filtering method

MI

Common Issues

Practical consideration

Summai

### Recall independence:

$$\frac{P(A,C)}{P(A)P(C)}=1$$

- $\,\blacksquare$  If LHS  $\sim$  1, attribute and class occur together as often as we would expect from random chance
- If LHS >> 1, attribute and class occur together much more often than randomly.
- (If LHS << 1, attribute and class are negatively correlated. More on this later.)



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Feature Selection

Wrappers Embedded

Embedded Filters

Filtering methods

PMI

 $\chi^2$ 

Common Issue

Practical consideration

Summary

Pointwise mutual information:

$$PMI(A, C) = \log_2 \frac{P(A, C)}{P(A)P(C)}$$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

Feature Selection

Wrappers Embedded

Embedded Filters

Filtering methods

PMI

 $\chi^2$ 

Common issues

Practical consideration

Summary

Pointwise mutual information:

$$PMI(A, C) = \log_2 \frac{P(A, C)}{P(A)P(C)}$$

Attributes with greatest PMI: best attributes



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnin

Eggtura Calagtic

Wrappers

Embedded

Filtering metho

PMI MI

Common Issue

Practical consideration

Summary

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	N
Ν	Ν	N

Calculate PMI of a<sub>1</sub>, a<sub>2</sub> with respect to c



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnin

#### Fasting Calcati

Wrappers Embedded

Filtoring mothe

### PMI

MI

Common Issue

Practical consideration

Summary

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	Ν
Ν	N	Ν

$$P(a_1) = \frac{2}{4}$$
;  $P(c) = \frac{2}{4}$ ;  $P(a_1, c) = \frac{2}{4}$ 



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Footure Coloctics

Wrappers

Embedded

....

#### Filtering metho

PMI

MI 2

00.....

Practical consideration

Summary

$$PMI(a_1, c) = log_2 \frac{\frac{1}{2}}{\frac{1}{2} \cdot \frac{1}{2}}$$
  
=  $log_2(2)$  =



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Heatures in Machine Learnin

#### Factoria Calcatia

Wrappers Embedded

Filters

Filtering metho

PMI MI

MI √²

Common Issue

Practical consideration

Summary

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	Ν
Ν	Ν	Ν

$$P(a_2) = \frac{2}{4}$$
;  $P(c) = \frac{2}{4}$ ;  $P(a_2, c) = \frac{1}{4}$ 



# Toy example, revisited

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Factions Calcution

Wrappers

Embedded

Filessins mades

PMI

МІ

00......

Practical consideration

Summary

$$PMI(a_2, c) = \log_2 \frac{\frac{1}{4}}{\frac{1}{2} \cdot \frac{1}{2}}$$
  
=  $\log_2(1)$  =



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

Wrappers Embedded

Filters

### Filtering methods

PMI MI

.....

Practical consideration

Summarı

What makes a single feature good?

Well correlated with class



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnii

#### Factions Calcutio

Wrappers Embedded

Filters

#### Filtering methods

PMI

MI

Common Is

Practical

Summary

What makes a single feature good?

- Well correlated with class
  - Knowing a lets us predict c with more confidence



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

#### Feature Selec

Embedded Filters

Filtering method

### PMI

PMI MI

Common Issue

Practical consideration

Summar

## What makes a single feature good?

- Well correlated with class
  - Knowing a lets us predict c with more confidence
- Reverse correlated with class
  - $\blacksquare$  Knowing  $\bar{a}$  lets us predict c with more confidence
  - Just as good



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection Wrappers

Embedded Filters

Filtering method

PMI MI

Common Issues

Practical consideration

Summary

### What makes a single feature good?

- Well correlated with class
  - Knowing a lets us predict c with more confidence
- Reverse correlated with class
  - $\blacksquare$  Knowing  $\bar{a}$  lets us predict c with more confidence
- Well correlated (or reverse correlated) with not class
  - Knowing a lets us predict c with more confidence
  - Usually not quite as good, but still useful



#### Lecture 18: **Feature Selection**

COMP90049 Technologies



Mutual information: combine each a, a, c, c PMI



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnir

Wrappers Embedded

Filters

Filtering method



Common Issue

Practical

Summar



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Feature Selecti

Embedded Eiltere

Filtering met

PMI



. . . .

consideration

Summary

$$\begin{array}{c|cc} & a & \bar{a} \\ \hline c & \sigma(a,c) & \sigma(\bar{a},c) \\ \bar{c} & \sigma(a,\bar{c}) & \sigma(\bar{a},\bar{c}) \\ \end{array}$$



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Feature Selecti

Wrappers Embedded Filters

Filters

PMI PMI



Common issu

Practical consideration

Summary

	а	ā	Total
С	$\sigma(a,c)$	$\sigma(\bar{a},c)$	$\sigma(c)$
$ar{c}$	$\sigma(a, \bar{c})$	$\sigma(\bar{\pmb{a}},\bar{\pmb{c}})$	$\sigma(\bar{c})$
Total	$\sigma(a)$	$\sigma(\bar{a})$	N



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learni

Feature Selection

Embedded Filters

Filtering method

PMI MI

Common

consideration

Summary

	а	ā	Total
С	$\sigma(a,c)$	$\sigma(\bar{a},c)$	$\sigma(c)$
$ar{c}$	$\sigma(a, \bar{c})$	$\sigma(\bar{\pmb{a}},\bar{\pmb{c}})$	$\sigma(ar{c})$
Total	$\sigma(a)$	$\sigma(\bar{a})$	N

$$P(a,c) = \frac{\sigma(a,c)}{N}$$
, etc.



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learn

Feature Selection

Wrappers Embedded

Filters

Filtering method

MI

 $\chi^2$ 

Described

consideration

Summar

Contingency tables for toy example:

$a_1$	$a_2$	С
Υ	Υ	Y
Υ	Ν	Υ
Ν	Υ	Ν
Ν	N	N

	$a_1$	a=Y	a = N	Total
	c =Y	2	0	2
	c = N	0	2	2
	Total	2	2	4
	$a_2$	a=Y	a = N	Total
_	<i>a</i> <sub>2</sub> <i>c</i> =Y	a =Y	a =N 1	Total 2
_		a=Y 1 1	a =N 1 1	
_	<i>c</i> =Y	a=Y 1 1 2	a =N 1 1 2	2



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learnin

Wrappers Embedded

Filters

Filtering metho

PMI

MI

Common Issues

Practical consideration

Summary

$$MI(A,C) = P(a,c)PMI(a,c) + P(\bar{a},c)PMI(\bar{a},c) + P(a,\bar{c})PMI(a,\bar{c}) + P(\bar{a},\bar{c})PMI(\bar{a},\bar{c})$$



#### Lecture 18: **Feature Selection**

COMP90049 Knowledge Technologies



$$MI(A, C) = P(a, c) \log_2 \frac{P(a, c)}{P(a)P(c)} + P(\bar{a}, c) \log_2 \frac{P(\bar{a}, c)}{P(\bar{a})P(c)} + P(\bar{a}, \bar{c}) \log_2 \frac{P(\bar{a}, \bar{c})}{P(a)P(\bar{c})} + P(\bar{a}, \bar{c}) \log_2 \frac{P(\bar{a}, \bar{c})}{P(\bar{a})P(\bar{c})}$$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

Feature Select

Wrappers Embedded

Filters

Filtering method

МІ

Common Iss

Practical consideration

Summary

Often written more compactly as:

$$MI(A, C) = \sum_{i \in \{a,\bar{a}\}} \sum_{j \in \{c,\bar{c}\}} P(i,j) \log_2 \frac{P(i,j)}{P(i)P(j)}$$

(This representation can be extended to different types of attributes more intuitively.)

Note that  $0 \log 0 \equiv 0$ .



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Feature Select

Embedded

Filters

Filtering metho

PMI

Commo

Practical

consideration

Summar

## Contingency table for toy example:

$a_1$	a=Y	a = N	Total
<i>c</i> =Y	2	0	2
c = N	0	2	2
Total	2	2	1

$$P(a, c) = \frac{2}{4}$$
;  $P(a) = \frac{2}{4}$ ;  $P(c) = \frac{2}{4}$   
 $P(\bar{a}, \bar{c}) = \frac{2}{4}$ ;  $P(\bar{a}) = \frac{2}{4}$ ;  $P(\bar{c}) = \frac{2}{4}$   
 $P(\bar{a}, c) = 0$ ;  $P(a, \bar{c}) = 0$ 

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	Ν
Ν	Ν	Ν



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

Wrappers Embedded

Filters

Filtering metho

PMI

MI

Common Issu

Practical consideration

Summary

$$\begin{aligned} MI(A_1,C) &= P(a_1,c)\log_2\frac{P(a_1,c)}{P(a_1)P(c)} + P(\bar{a}_1,c)\log_2\frac{P(\bar{a}_1,c)}{P(\bar{a}_1)P(c)} + \\ &P(a_1,\bar{c})\log_2\frac{P(a_1,\bar{c})}{P(a_1)P(\bar{c})} + P(\bar{a}_1,\bar{c})\log_2\frac{P(\bar{a}_1,\bar{c})}{P(\bar{a}_1)P(\bar{c})} \\ &= \frac{1}{2}\log_2\frac{\frac{1}{2}}{\frac{1}{2}\frac{1}{2}} + 0\log_2\frac{0}{\frac{1}{2}\frac{1}{2}} + 0\log_2\frac{0}{\frac{1}{2}\frac{1}{2}} + \frac{1}{2}\log_2\frac{\frac{1}{2}}{\frac{1}{2}\frac{1}{2}} \\ &= \frac{1}{2}(1) + 0 + 0 + \frac{1}{2}(1) = 1 \end{aligned}$$



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Feature Selec

Wrappers Embedded Filters

Filtering method

PMI MI

Common Is

Practical

consideration

Summar

## Contingency table for toy example:

$a_2$	a =Y	a = N	Total
c=Y	1	1	2
c = N	1	1	2
Total	2	2	4

$$P(a,c) = \frac{1}{4}$$
;  $P(a) = \frac{2}{4}$ ;  $P(c) = \frac{2}{4}$   
 $P(\bar{a},\bar{c}) = \frac{1}{4}$ ;  $P(\bar{a}) = \frac{2}{4}$ ;  $P(\bar{c}) = \frac{2}{4}$   
 $P(\bar{a},c) = \frac{1}{4}$ ;  $P(a,\bar{c}) = \frac{1}{4}$ 

$a_1$	$a_2$	С
Υ	Υ	Υ
Υ	Ν	Υ
Ν	Υ	Ν
Ν	Ν	Ν



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

Embedded

Filters

Filtering method

PMI

МІ

Common Iss

Practical consideration

Summary

$$\begin{split} \textit{MI}(\textit{A}_{2},\textit{C}) &= P(\textit{a}_{2},\textit{c})\log_{2}\frac{P(\textit{a}_{2},\textit{c})}{P(\textit{a}_{2})P(\textit{c})} + P(\bar{\textit{a}}_{2},\textit{c})\log_{2}\frac{P(\bar{\textit{a}}_{2},\textit{c})}{P(\bar{\textit{a}}_{2})P(\textit{c})} + \\ &P(\textit{a}_{2},\bar{\textit{c}})\log_{2}\frac{P(\textit{a}_{2},\bar{\textit{c}})}{P(\textit{a}_{2})P(\bar{\textit{c}})} + P(\bar{\textit{a}}_{2},\bar{\textit{c}})\log_{2}\frac{P(\bar{\textit{a}}_{2},\bar{\textit{c}})}{P(\bar{\textit{a}}_{2})P(\bar{\textit{c}})} \\ &= \frac{1}{4}\log_{2}\frac{\frac{1}{4}}{\frac{1}{2}\frac{1}{2}} + \frac{1}{4}\log_{2}\frac{\frac{1}{4}}{\frac{1}{2}\frac{1}{2}} + \frac{1}{4}\log_{2}\frac{\frac{1}{4}}{\frac{1}{2}\frac{1}{2}} + \frac{1}{4}\log_{2}\frac{\frac{1}{4}}{\frac{1}{2}\frac{1}{2}} \\ &= \frac{1}{4}(0) + \frac{1}{4}(0) + \frac{1}{4}(0) = 0 \end{split}$$



# Chi-square

Lecture 18: Feature Selection COMP90049 Knowledge

Technologies

Features in
Machine Learning

Feature Selection

Embedded Filters

Filtering method

MI

Common Issue

Practical considerations

Summs

Check the value we actually observed O(W) with the expected value E(W):

- If the observed value is much greater than the expected value, a occurs more often with c than we would expect at random predictive
- If the observed value is much lesser than the expected value, a occurs less often with c than we would expect at random predictive
- If the observed value is close to the expected value, a occurs as often with c as we would expect randomly — not predictive

Similarly with X, Y, Z



# Chi-square

#### Lecture 18: **Feature Selection**

COMP90049 Knowledge Technologies

Actual calculation (written more compactly):

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}$$

(*i* sums over rows and *j* sums over columns.)

In practice, there are simpler ways to calculate this for  $2 \times 2$  contingency tables.



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

#### Feature Selection

Embedded Filters

Filtering method

PMI

#### **Common Issues**

Practical consideration

Summary

So far, we've only looked at binary (Y/N) attributes:

- Nominal attributes
- Continuous attributes
- Ordinal attributes



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

Embedded Eiltere

Filtering method

PMI

MI

Common Issues

Practical consideration

Summar

Nominal attributes (e.g. Outlook={sunny, overcast, rainy}). Two common strategies:

- Treat as multiple binary attributes:
  - e.g. sunny=Y, overcast=N, rainy=N, etc.
  - Can just use the formulae as given
  - Results often difficult to interpret
    - For example, Outlook=sunny is useful, but Outlook=overcast and Outlook=rainy are not useful... Should we use Outlook?
- Modify contigency tables (and formulae)



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection
Wrappers

Embedded Filters

Filtering method

**Common Issues** 

Practical consideration

Summar

## Modified contingency table:

$$\begin{array}{c|cccc} \textbf{0} & \textbf{s} & \textbf{o} & \textbf{r} \\ \hline c = \textbf{Y} & \textbf{U} & \textbf{V} & \textbf{W} \\ c = \textbf{N} & \textbf{X} & \textbf{Y} & \textbf{Z} \\ \end{array}$$

Modified MI:

$$\begin{split} \mathit{MI}(O,C) &= \sum_{i \in \{s,o,r\}} \sum_{j \in \{c,\bar{c}\}} P(i,j) \log_2 \frac{P(i,j)}{P(i)P(j)} \\ &= P(s,c) \log_2 \frac{P(s,c)}{P(s)P(c)} + P(s,\bar{c}) \log_2 \frac{P(s,\bar{c})}{P(s)P(\bar{c})} + \\ &= P(o,c) \log_2 \frac{P(o,c)}{P(o)P(c)} + P(o,\bar{c}) \log_2 \frac{P(o,\bar{c})}{P(o)P(\bar{c})} + \\ &= P(r,c) \log_2 \frac{P(r,c)}{P(r)P(c)} + P(r,\bar{c}) \log_2 \frac{P(r,\bar{c})}{P(r)P(\bar{c})} \end{split}$$

Biased towards attributes with many values. (Why?)





### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection Wrappers

Embedded Filters

Filtering method

 $\chi^2$ 

Common Issues

Practical consideration

Summa

### Continuous attributes:

- Usually dealt with by estimating probability based on a Gaussian (normal) distribution
- With a large number of values, most random variables are normally distributed due to the Central Limit Theorem
- For small data sets or pathological features, we typically need to use messy binomial/multinomial distributions

All of this is (unsurprisingly) beyond the scope of this subject



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learr

Feature Selection Wrappers

Embedded Filters

Filtering method

Common Issues

Practical consideration

Summar

Ordinal attributes (e.g. low, med, high or 1,2,3,4). Three possibilities, roughly in order of popularity:

- Treat as binary
  - Particularly appropriate for frequency counts where events are low-frequency (e.g. words in tweets)
- Treat as nominal (i.e. throw away ordering)
- Treat as continuous
  - The fact that we haven't seen any intermediate values is usually not important
  - Does have all of the technical downsides of continuous attributes, however



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

Embedded Filters

Filtering method

PMI

Common Issues

Practical consideration

.....

So far, we've only looked at binary (Y/N) classification tasks.

What makes a single feature good?

- Highly correlated with class
- Highly reverse correlated with class
- Highly correlated (or reverse correlated) with not class
- ... What if there are many classes?

Multiclass (e.g., RightTroll, LeftTroll, and Other) classification tasks are usually much more difficult.



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Fasture Calcatia

Wrappers Embedded

Filters

# Filtering method

PMI MI

#### Common Issues

Practical consideration

Summary

What makes a feature bad?

Irrelevant



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Fasture Calcatia

Embedded

Filtering method

PMI

MI

#### **Common Issues**

Practical consideration

Summary

### What makes a feature **bad**?

- Irrelevant
- Correlated with other features



#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnir

#### Feature Selection

Embedded Filters

Filtering method

PMI

 $\chi^2$ 

### Common Issues

Practical consideration

Summary

### What makes a feature bad?

- Irrelevant
- Correlated with other features
- Good at only predicting one class (but is this truly bad?)



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selection Wrappers Embedded

Filtering methods

MI 2

Common Issues

Practical consideration

Summary

### Consider multi-class problem over R, L, O:

- PMI, MI,  $\chi^2$  are all calculated *per-class*
- (Some other feature selection metrics, e.g. Information Gain, work for all classes at once)
- Need to make a point of selecting (hopefully uncorrelated) features for each class to give our classifier the best chance of predicting everything correctly.



### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Wrappers Embedded

Filtering metho

PMI MI

Common Issues

Practical consideration

Summar

## Actual example (MI):

R	L	0
breaking	the	the
maga	and	to
rtamerica	a	is
retweet	black	a
beeth	i	and
bbsp	is	this
tcot	news	rt
obama	to	of
hillary	blacklivesmatter	breaking
pjnet	this	on



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnin

Wrappers Embedded

Filtering met

PMI MI

 $\chi^2$ 

Common Issues

Practical consideration

Summar

# Actual example ( $\chi^2$ ):

R	L	0
breaking	the	the
maga	and	to
rtamerica	a	is
retweet	black	a
beeth	i	and
tcot	is	of
bbsp	blacklivesmatter	this
obama	this	rt
hillary	to	on
rt	amp	dallas



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learni

Feature Selection

Embedded

Filters

PMI

MI √²

Common Issues

Practical consideration

Summar

## Actual example (MI):

	В	Н	Se	SD	W
_	boston	houston	seattle	diego	dc
	diego	diego	diego	san	diego
	san	jupdicom	wa	chargers	san
	httpbitlyczmk	tx	san	sd	obama
	ma	san	cheezburger	sdut	health
	redsox	httpbitlycdqk	boston	seattle	washington
	seattle	seattle	bellevue	sandiego	bill



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learn

Feature Selection

Wrappers Embedded

Filters

Filtering methods

MI √²

Common Issues

Practical consideration

Summai

### Intuitive features:

В	Н	Se	SD	W
boston	houston	seattle	diego	dc
diego	diego	diego	san	diego
san	jupdicom	wa	chargers	san
httpbitlyczmk	tx	san	sd	obama
ma	san	cheezburger	sdut	health
redsox	httpbitlycdqk	boston	seattle	washington
		le ellerere		1. 111
seattle	seattle	bellevue	sandiego	bill



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learn

Feature Selecti

Wrappers Embedded

Filters

Filtering method

MI ײ

Common Issues

Practical consideration

Summar

Features for predicting not class (MI):

В	Н	Se	SD	W
boston	houston	seattle	diego	dc
diego	diego	diego	san	diego
san	jupdicom	wa	chargers	san
httpbitlycz	mk tx	san	sd	obama
ma	san	cheezburger	sdut	health
redsox	httpbitlycdqk	boston	seattle	washington
seattle	seattle	bellevue	sandiego	bill



Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learnii

Feature Selection

Embedded

Filters

Filtering methods

MI 2

 $\chi^{\epsilon}$  Common Issues

00.....

consideration

Summar

### Unintuitive features:

В	Н	Se	SD	W
boston	houston	seattle	diego	dc
diego	diego	diego	san	diego
san	jupdicom	wa	chargers	san
httpbitlyczmk	tx	san	sd	obama
ma	san	cheezburger	sdut	health
redsox	httpbitlycdqk	boston	seattle	washington
seattle	seattle	bellevue	sandiego	bill



# What's going on with MI?

### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection Wrappers Embedded Filters

Filtering method PMI MI

Common Issue

Practical considerations

Summary

Mutual Information is biased toward rare, uninformative features

- All probabilities: no notion of the raw frequency of events
- If a feature is seen rarely, but always with a given class, it will be seen as "good"
- For example: httpbitlyczmk occurs 447 times out of 750K instances, but often with B. Is this meaningful?
- Best features in the Twitter dataset only had MI of about 0.1 bits; 100<sup>th</sup> best for a given class had MI of about 0.0001 bits



# So... Give up on feature selection then?

### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learning

Feature Selection

Embedded Filters

Filtering method

PMI

Common Is:

Practical considerations

Summary

### No way!

- Even marginally relevant features usually a vast improvement on an unfiltered data set
- Some models need feature selection
  - k-Nearest Neighbour, hugely
  - Naive Bayes, Decision Trees, andSVM to a lesser extent
- Machine learning experts (us!) need to think about the data!



# Summary

#### Lecture 18: Feature Selection

COMP90049 Knowledge Technologies

Machine Learni

Feature Selecti

Wrappers Embedded Filters

Filtering method

 $\chi^2$ 

Common issues

Practical consideration

Summary

- Wrappers vs. Embedded methods vs. Filters
- Popular filters: PMI, MI,  $\chi^2$ , how should we use them and what are the results going to look like
- Importance of feature selection for different methods (even though it often isn't the solution we were hoping for)



# References

#### Lecture 18: Feature Selection COMP90049 Knowledge Technologies

Features in Machine Learning

Feature Selection Wrappers Embedded

Filtering method

PMI MI

Common Issues

Practical consideration

Summary

Guyon, Isabelle, and Andre Elisseeff. 2003. An introduction to variable and feature selection. *The Journal of Machine Learning Research*. Vol 3, 1157–1182.

John, George, Ron Kohavi, and Karl Pfleger. 1994. Irrelevant features and the subset selection problem. In *Proceedings of the 11th International Conference on Machine Learning*, 121–9.

Tan, Pang-Ning, Michael Steinbach, and Vipin Kumar. 2006. *Introduction to Data Mining.* Addison Wesley.

Witten, Ian, and Eibe Frank. 2005. *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. San Francisco, USA: Morgan Kaufmann.

Yang, Yiming and Jan Pedersen. 1997. A Comparative Study on Feature Selection in Text Categorization. In *Proceedings of the Fourteenth International Conference on Machine Learning*, 412–20.