

# Text Processing - Basics

Advanced Social Computing

Department of Computer Science  
University of Massachusetts, Lowell  
Spring 2020

Hadi Amiri  
[hadi@cs.uml.edu](mailto:hadi@cs.uml.edu)



# Announcement

- **SP1 out**
  - Due Date: 3/25, 3:30 PM

# Lecture Topics

- Text Data
- Learning word vectors
  - Word2vec
  - Glove
- Evaluating word vectors
- Retrofitting word vectors

# Text Processing!

- Tokenization
- Normalization
  - urls, hashtags, punctuations, numbers, dates, cases, stop words, etc.
  - Spell correction
- Morphological analysis
  - Stemming, lemmatization, etc.
- Syntactic analysis
- Semantic analysis
- Discourse analysis

# Text Processing!

- Tokenization
- Normalization
- Morphology
- Syntax
- Semantic
- Discourse

**Text cleaning is a very important first step!  
But there is no general rule.**

- Is it safe to remove punctuations or stop words from text?
  - “switching from Verizon” vs. “switching to Verizon.”
- Or convert all characters to lowercase?
  - “Bush” vs. “bush.”
- Or remove all numbers?
  - “7 yrs old” vs. “70 yrs old.”

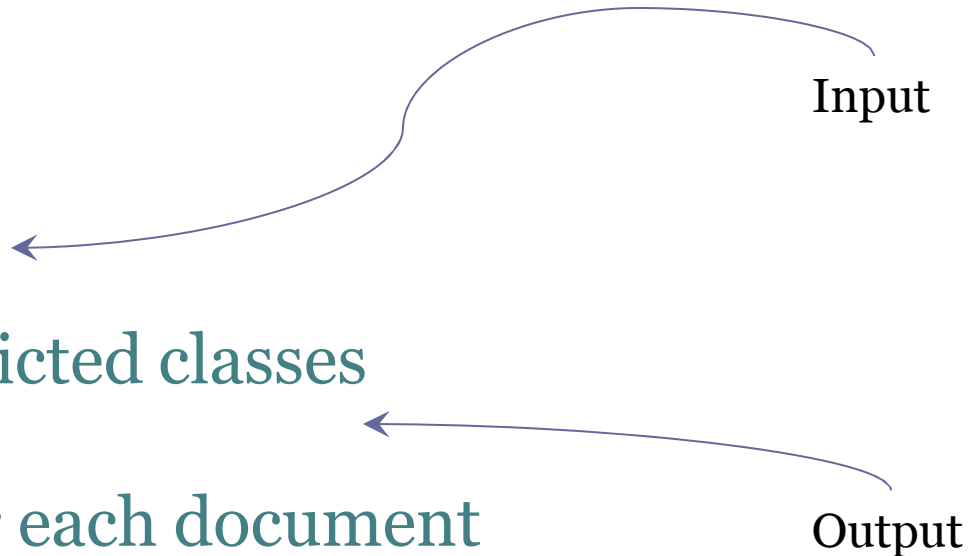
- **General rule:** Use the exact same cleaning technique for all competing models.

# Many Interesting Applications

- Search
- Information Extraction
- Question Answering
- Machine Translation
- Summarization
- Dialogue Systems
- Text Classification
  - **Emails:** spam, not-spam.
  - **News articles:** business, health, sports, tech, etc.
  - **Reviews:** positive, negative, neutral.
  - **Word pairs:** synonyms or not.
  - **Essays as:** A, B, C, D, or F
  - **Etc.**

# Text Classification

- Let's say we have:
  - A set of documents
    - $\mathbf{X} = \{x_1, \dots, x_n\}$
  - A set of labels or predicted classes
    - $\mathbf{Y} = \{\text{Class-1}, \dots, \text{Class-k}\}$
  - We know the label for each document
    - $(x_1, y_1), \dots, (x_n, y_n)$
  - We aim to learn a function  $f$  (classifier) that can map inputs to their corresponding outputs
    - $f: \mathbf{X} \rightarrow \mathbf{Y}$



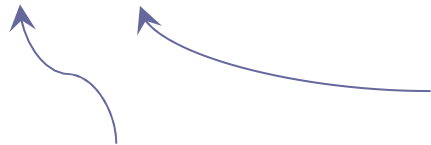
# Text Classification- Cnt.

- $\mathbf{X} = \{$  i love verizon's coverage ,  
 actually t-mobile has great deals,  
 i hate t-mobile! One more bill!!,  
 i cant take it anymore! hate verizon  $\}$

- $\mathbf{Y} = \{+1, -1\}$

**positive**

**negative**





# Text Classification- Cnt.

- $\mathbf{X} = \{$  i love verizon's coverage ,  
actually t-mobile has great deals,  
i hate t-mobile! One more bill!!,  
i cant take it anymore! hate verizon  $\}$
- $\mathbf{Y} = \{+1, -1\}$
- $\{(x_1, y_1), (x_2, y_2), \dots, (x_4, y_4)\} =$   
 $\{($  i love verizon's coverage,  $+1)$ ,  
 $($  actually t-mobile has great deals,  $+1)$ ,  
 $($  i hate t-mobile! One more bill!!,  $-1)$ ,  
 $($  i cant take it anymore! hate verizon,  $-1)$   $\}$

# Text Classification- Cnt.

- $\mathbf{X} = \{$  i love verizon's coverage ,  
actually t-mobile has great deals,  
i hate t-mobile! One more bill!!,  
i cant take it anymore! hate verizon  $\}$

- $\mathbf{Y} = \{+1, -1\}$

- $f(\text{i love verizon's coverage}) = +1$

- $f(\text{actually t-mobile has great deals}) = +1$

- $f(\text{i hate t-mobile! One more bill!!}) = -1$

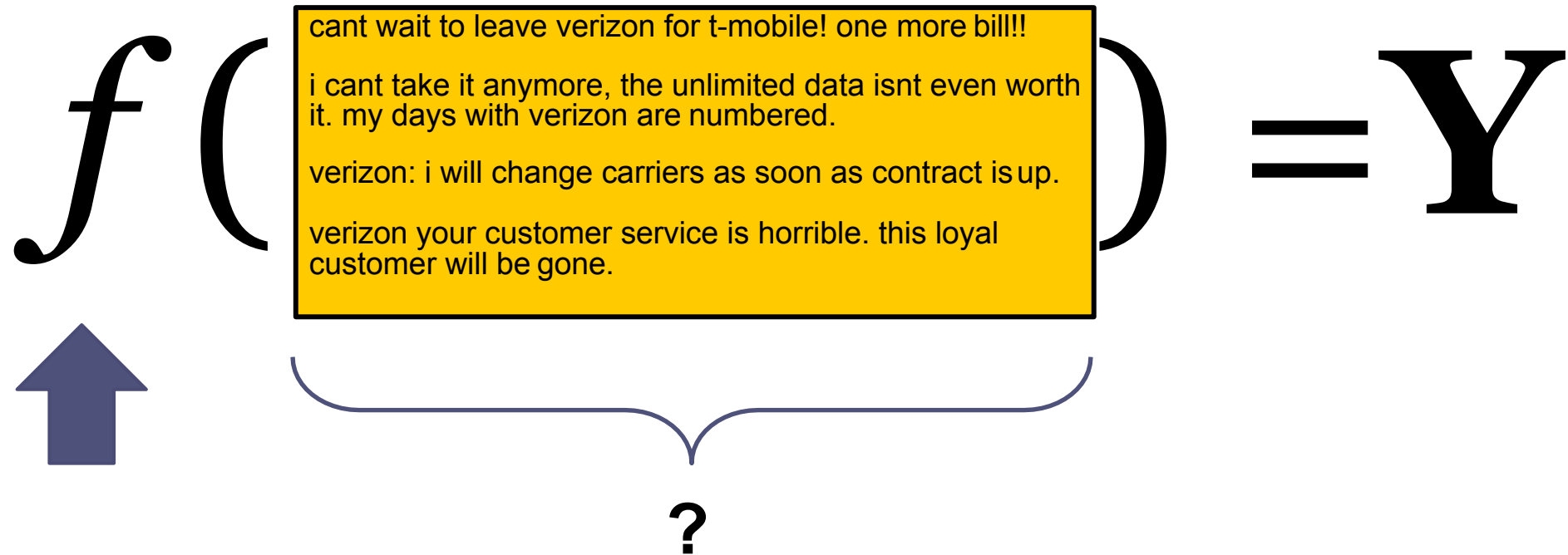
- $f(\text{i cant take it anymore! hate verizon}) = -1$

**Classification:** the output variable takes class labels, i.e.  $Y = \{-1, +1\}$

**Regression:** the output variable takes continuous values, i.e.  $Y = [-1, +1]$ .

Why do we need  
to learn  $f$ ?

# Text Classification- Cnt.



How can we determine  $f$ ?

# How can we determine $f$

- Given  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ , we aim to find  $f(\cdot)$ !
- An ideal  $f(\cdot)$  is a function such that
  - $f(x_i) = y_i$  for all  $i$
  - Hard to find, why?
  - We just expect  $f(x_i)$  to be very close to  $y_i$ .

$y$	$f(x)$	$error = (y - f(x))^2$
+1	+1	0
-1	-1	0
+1	-1	4
-1	+1	4

Loss function

$$l(y_i, f(x_i)) = (y_i - f(x_i))^2$$

$$\sum_i l(y_i, f(x_i)) = \sum_i (y_i - f(x_i))^2$$

zero error when prediction and actual label are the same,  
Non-zero, otherwise!

# How can we determine $f$ -Cnt.

- Given  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  we aim to find a function  $f(\cdot)$  that minimizes the error

$$L(x, y) = \sum_i^n l(y_i, f(x_i))$$

- Three popular loss functions
  - Squared loss (linear classifier)
  - Hinge loss (the SVMs),
  - Logistic loss (logistic classifier)

# How can we determine $f$ - Cnt.

- Three popular loss functions  $L(x, y) = \sum_i^n l(y_i, f(x_i))$ 
  - Squared loss (linear classifier)

$$l(y, f(x)) = (y - f(x))^2$$

- Hinge loss (the SVMs)

$$l(y, f(x)) = \max(0, 1 - y \cdot f(x))$$

- Logistic loss (logistic classifier)

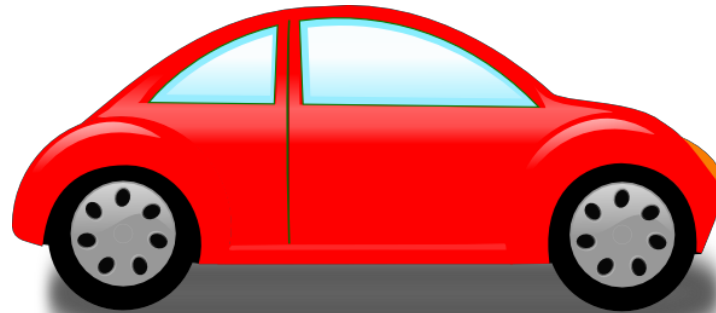
$$l(y, f(x)) = \log(1 + \exp(-y \cdot f(x)))$$

# Text Representation

- What is a good way to represent the input text?

# Text Representation- Cnt.

- Features
  - How to classify objects such as **People** and **Cars**?

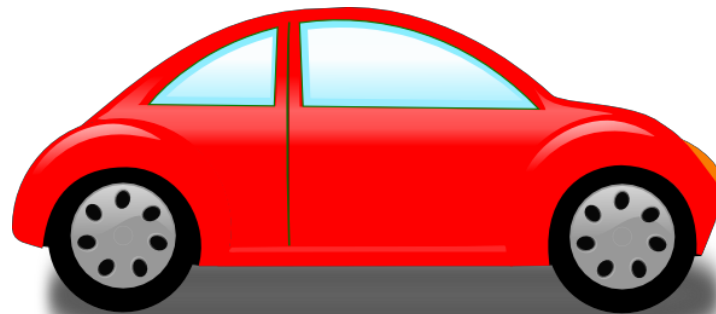


- We use *features / characteristics* of those objects!



# Text Representation- Cnt.

- Knowledge about features that make good predictors of class membership!
  - having wheels or not distinguishes people from cars, but doesn't distinguish cars from planes.



# Text Representation- Cnt.

- $\mathbf{X} = \{$  i would love verizon coverage,  
i hate verizon one more bill,  
i hate verizon  $\}$

- $\mathbf{Y} = \{+1, -1\}$

Bag of Word representation

- Features = [i, would, love, verizon, coverage, hate, one, more, bill]

# Text Representation- Cnt.

- $\mathbf{X} = \{$  i would love verizon coverage,  
i hate verizon one more bill,  
i hate verizon  $\}$

- $\mathbf{Y} = \{+1, -1\}$

Bag of Word representation

- Features = [i, would, love, verizon, coverage, hate, one, more, bill]

Feature Weights

	i	would	love	verizon	coverage	hate	one	more	bill
$x_1$	1	1	1	1	1	0	0	0	0
$x_2$	1	0	0	1	0	1	1	1	1
$x_3$	1	0	0	1	0	1	0	0	0

# Text Representation- Cnt.

- Bag of Word representation
- $\mathbf{X} = \{$  i would love verizon coverage,  
i hate verizon one more bill,  
i hate verizon  $\}$
- $\mathbf{Y} = \{+1, -1\}$
- Features = {would, love, hate}

Sentiment Words



	would	love	hate
$x_1$	1	1	0
$x_2$	0	0	1
$x_3$	0	0	1

# Text Representation- Cnt.

- Other ways of representation?
- Other ways to set weights?
- How to encode semantics?
  - Suggest & recommend
  - Pretty & beautiful
  - Etc.

# Vowpal Wabbit (VW)

- Vowpal Wabbit:
  - Fast learning
  - Simplicity
  - Namespace definition
    - Easy Ablation Analysis



<http://hunch.net/~vw/>

Label [Importance] [Base] ['Tag] | Namespace Feature ... | Namespace Feature ...

Namespace = A letter like 'a', 'b', 'c', ...

Feature = String[:Float]

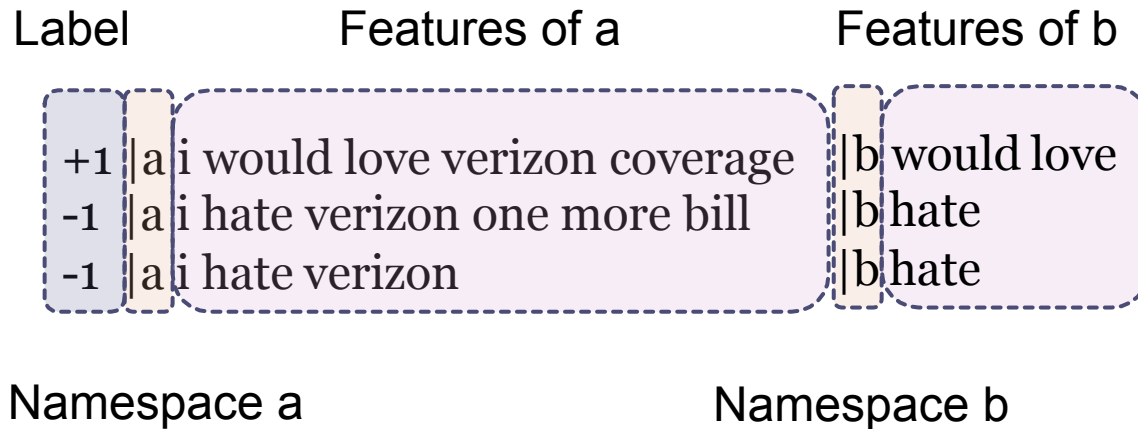
# Vowpal Wabbit (VW)

- Vowpal Wabbit:
  - Fast learning
  - Simplicity
  - Namespace definition
    - Easy Ablation Analysis

+1	a i would love verizon coverage	b would love
-1	a i hate verizon one more bill	b hate
-1	a i hate verizon	b hate

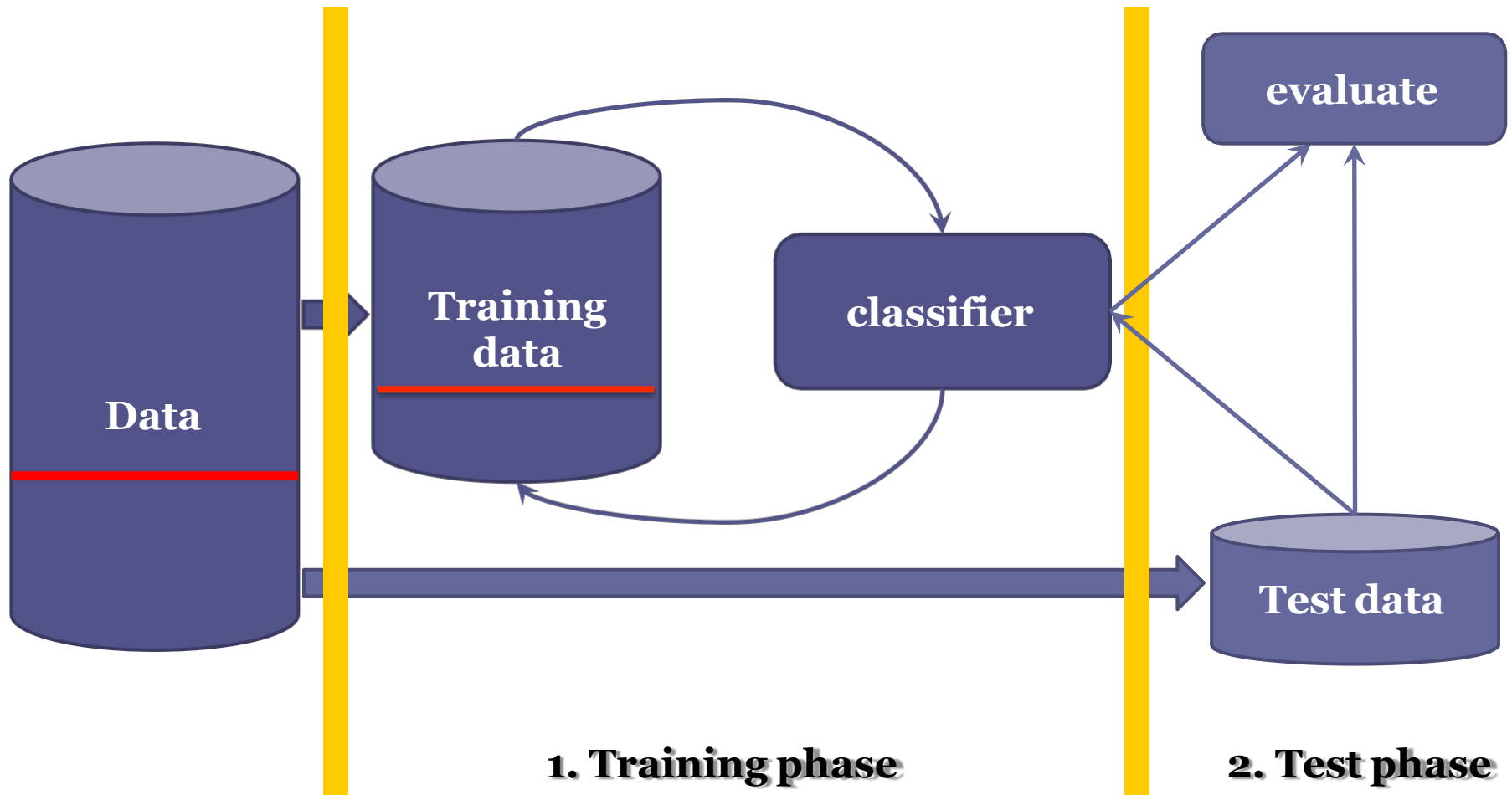
# Vowpal Wabbit (VW)

- Vowpal Wabbit:
  - Fast learning
  - Simplicity
  - Namespace definition
    - Easy Ablation Analysis



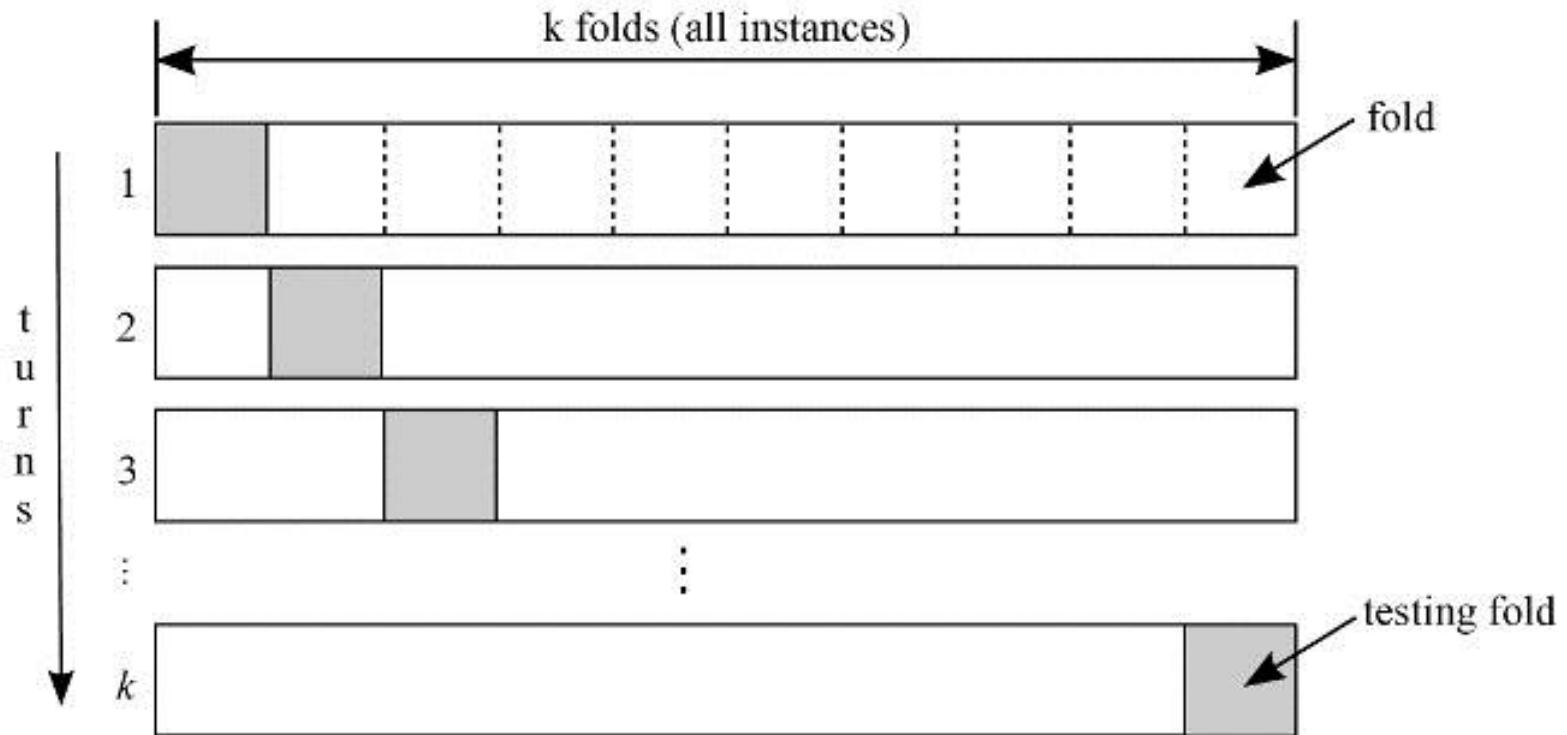


# Test and Training Data



# Test and Training Data- Cnt.

- How to create test and training data?
  - Use k-fold cross validation,  $k=3$  or  $5$



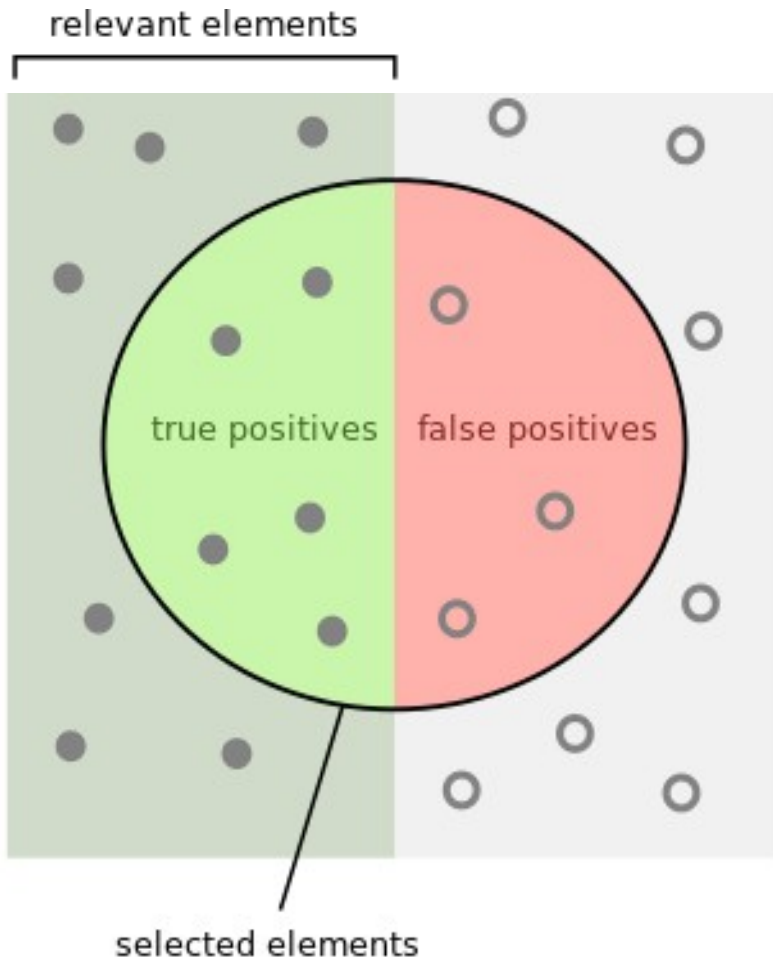
# Evaluation

- Commonly-used evaluation metrics


Scoring	Function	Comment
<b>Classification</b>		
'accuracy'	<code>metrics.accuracy_score</code>	
'balanced_accuracy'	<code>metrics.balanced_accuracy_score</code>	
'average_precision'	<code>metrics.average_precision_score</code>	
'neg_brier_score'	<code>metrics.brier_score_loss</code>	
'f1'	<code>metrics.f1_score</code>	for binary targets
'f1_micro'	<code>metrics.f1_score</code>	micro-averaged
'f1_macro'	<code>metrics.f1_score</code>	macro-averaged
'f1_weighted'	<code>metrics.f1_score</code>	weighted average
'f1_samples'	<code>metrics.f1_score</code>	by multilabel sample
'neg_log_loss'	<code>metrics.log_loss</code>	requires <code>predict_proba</code> support
'precision' etc.	<code>metrics.precision_score</code>	suffixes apply as with 'f1'
'recall' etc.	<code>metrics.recall_score</code>	suffixes apply as with 'f1'
'jaccard' etc.	<code>metrics.jaccard_score</code>	suffixes apply as with 'f1'
'roc_auc'	<code>metrics.roc_auc_score</code>	
'roc_auc_ovr'	<code>metrics.roc_auc_score</code>	
'roc_auc_ovo'	<code>metrics.roc_auc_score</code>	
'roc_auc_ovr_weighted'	<code>metrics.roc_auc_score</code>	
'roc_auc_ovo_weighted'	<code>metrics.roc_auc_score</code>	
<b>Clustering</b>		
'adjusted_mutual_info_score'	<code>metrics.adjusted_mutual_info_score</code>	
'adjusted_rand_score'	<code>metrics.adjusted_rand_score</code>	
'completeness_score'	<code>metrics.completeness_score</code>	
'fowlkes_mallows_score'	<code>metrics.fowlkes_mallows_score</code>	
'homogeneity_score'	<code>metrics.homogeneity_score</code>	
'mutual_info_score'	<code>metrics.mutual_info_score</code>	
'normalized_mutual_info_score'	<code>metrics.normalized_mutual_info_score</code>	
'v_measure_score'	<code>metrics.v_measure_score</code>	
<b>Regression</b>		
'explained_variance'	<code>metrics.explained_variance_score</code>	
'max_error'	<code>metrics.max_error</code>	
'neg_mean_absolute_error'	<code>metrics.mean_absolute_error</code>	
'neg_mean_squared_error'	<code>metrics.mean_squared_error</code>	
'neg_root_mean_squared_error'	<code>metrics.mean_squared_error</code>	
'neg_mean_squared_log_error'	<code>metrics.mean_squared_log_error</code>	
'neg_median_absolute_error'	<code>metrics.median_absolute_error</code>	

# Evaluation – Cnt.


- Precision, Recall, F1-Score



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$


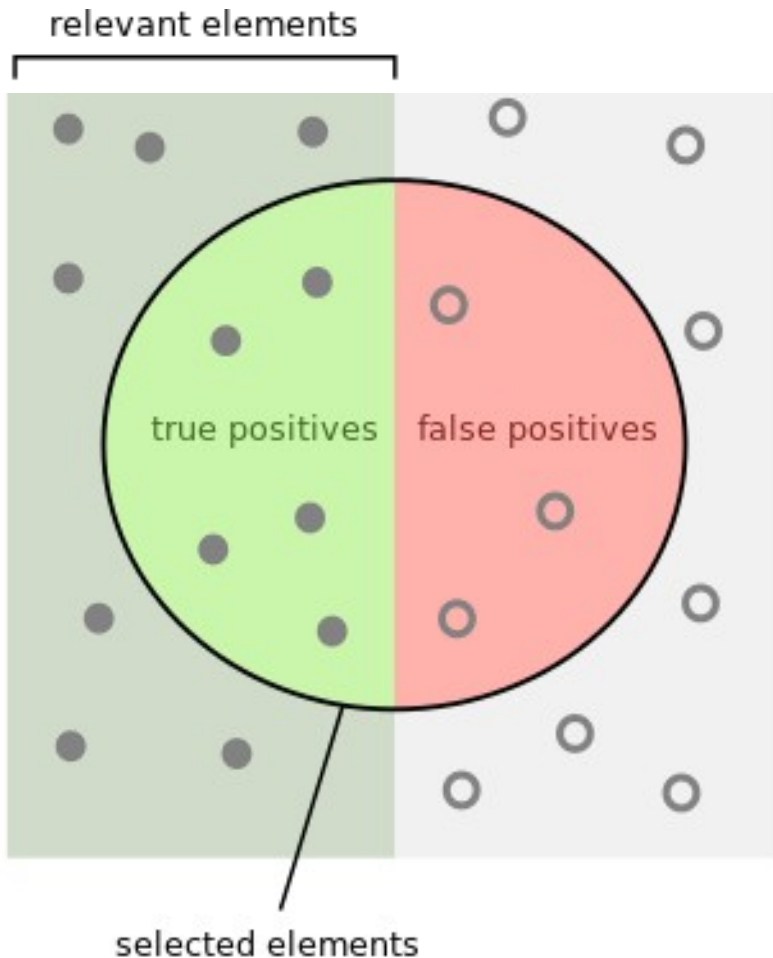
How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$



F1 is the harmonic mean of precision and recall

# Evaluation – Cnt.


- Precision, Recall, F1-Score



How many selected items are relevant?

$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$


How many relevant items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$


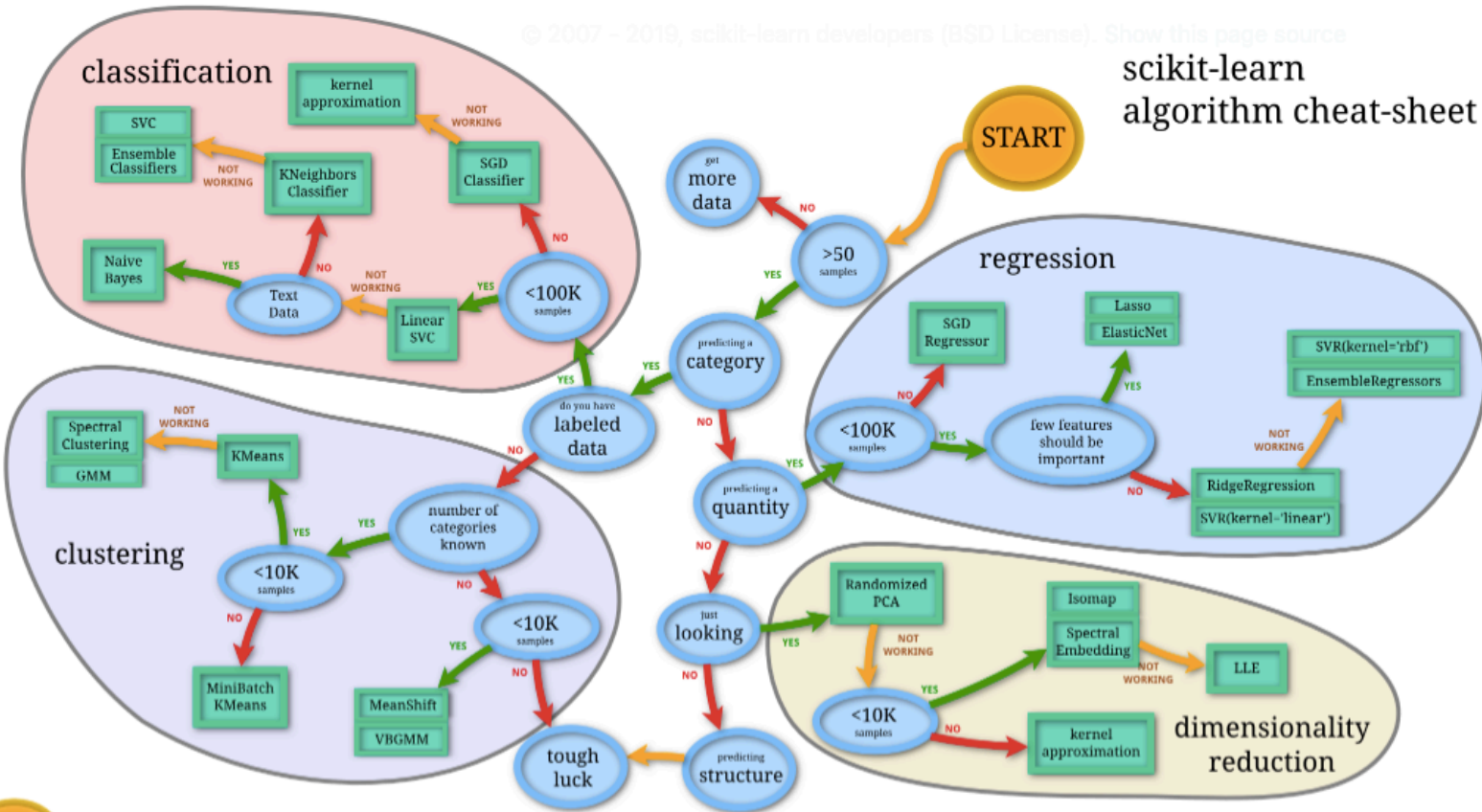
F1 is the harmonic mean of precision and recall

In which applications precision or recall is more important than the other?

# Quick Reference

© 2007 - 2019, scikit-learn developers (BSD License). Show this page source

scikit-learn  
algorithm cheat-sheet



# Questions?