

Text Classification

COMP90042 Lecture 7



THE UNIVERSITY OF
MELBOURNE

Outline

- Fundamentals of classification
- Text classification tasks
- Algorithms for classification
- Evaluation

Classification

- Input
 - * A document d
 - Often represented as a vector of *features*
 - * A fixed output set of classes $C = \{c_1, c_2, \dots, c_k\}$
 - Categorical, not continuous (regression) or ordinal (ranking)
- Output
 - * A predicted class $c \in C$

Text classification tasks

- Obviously more than can be enumerated here
- But let's briefly look at some major examples:
 - * Topic classification
 - * Sentiment analysis
 - * Authorship attribution
 - * Native-language identification
 - * Automatic fact-checking
- Not necessarily a full-sized "text"
 - * E.g. sentence or tweet-level polarity classification

Topic classification

- Motivation: library science, information retrieval
- Classes: Topic categories, e.g. “jobs”, “anxiety disorders”
- Features
 - * *Unigram* bag of words (BOW), with stop-words removed
 - * Longer *n-grams* (*bigrams, trigrams*) for phrases
- Examples of corpora
 - * Reuters news corpus (RCV1, see NLTK sample)
 - * Pubmed abstracts
 - * Tweets with hashtags

Topic classification Example

Is the topic of this text from the Reuters news corpus acquisitions or earnings?

LIEBERT CORP APPROVES MERGER

Liebert Corp said its shareholders approved the merger of a wholly-owned subsidiary of Emerson Electric Co. Under the terms of the merger, each Liebert shareholder will receive .3322 shares of Emerson stock for each Liebert share.

ANSWER: ACQUISITIONS

Sentiment Analysis

- Motivation: opinion mining, business analytics
- Classes: Positive/Negative/(Neutral)
- Features
 - * *N*-grams
 - * Polarity lexicons
- Examples of corpora
 - * Polarity movie review dataset (in NLTK)
 - * SEMEVAL Twitter polarity datasets

Sentiment analysis example

What is the polarity of this tweet from the SEMEVAL dataset?

anyone having problems with Windows 10? may be coincidental but since i downloaded, my WiFi keeps dropping out. Itunes had a malfunction

ANSWER: NEGATIVE

Authorship attribution

- Motivation: forensic linguistics, plagiarism detection
- Classes: Authors (e.g. Shakespeare)
- Features
 - * Frequency of function words
 - * Character n -grams
 - * Discourse structure
- Examples of corpora
 - * Project Gutenberg corpus (see NLTK sample)

Author attribution example

Which famous novelist wrote this text from Project Gutenberg?

Mr. Dashwood's disappointment was, at first, severe; but his temper was cheerful and sanguine; and he might reasonably hope to live many years, and by living economically, lay by a considerable sum from the produce of an estate already large, and capable of almost immediate improvement. But the fortune, which had been so tardy in coming, was his only one twelvemonth. He survived his uncle no longer; and ten thousand pounds, including the late legacies, was all that remained for his widow and daughters.

ANSWER: JANE AUSTEN

Native-Language Identification

- Motivation: forensic linguistics, educational applications
- Classes: first language of author (e.g. Chinese)
- Features
 - * Word N -grams
 - * Syntactic patterns (POS, parse trees)
 - * Phonological features
- Examples of corpora
 - * TOEFL/IELTS essay corpora

Native-Language Identification

What is the native language of the writer of this text?

Now a festival of my university is being held, and my club is joining it by offering a target practice game using bows and arrows of archery. I'm a manager of the attraction, so I have worked to make it succeed. I found it puzzled to manage a event or a program efficiently without generating a free rider. The event is not free, so we earn a lot of money.

ANSWER: JAPANESE

Automatic fact-checking

- Motivation: social media, journalism (fake news)
- Classes: True/False/(Can't be sure)
- Features
 - * N-grams
 - * Non-text metadata
 - * ??? (very recent task)
- Examples of corpora
 - * Emergent, LIAR: political statements
 - * FEVER

Automatic fact-checking

Is this statement true or false?

Austin is burdened by the fastest-growing tax increases of any major city in the nation.

ANSWER:FALSE

Building a Text classifier

1. Identify a task of interest
2. Collect an appropriate corpus
3. Carry out annotation
4. Select features
5. Choose a machine learning algorithm
6. Tune hyperparameters using held-out development data
7. Repeat earlier steps as needed
8. Train final model
9. Evaluate model on held-out test data

Choosing a classification algorithm

- Bias vs. Variance
- Feature independence
- Feature scaling
- Complexity
- Speed

Naïve Bayes

- Finds the class with the highest likelihood under Bayes law

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

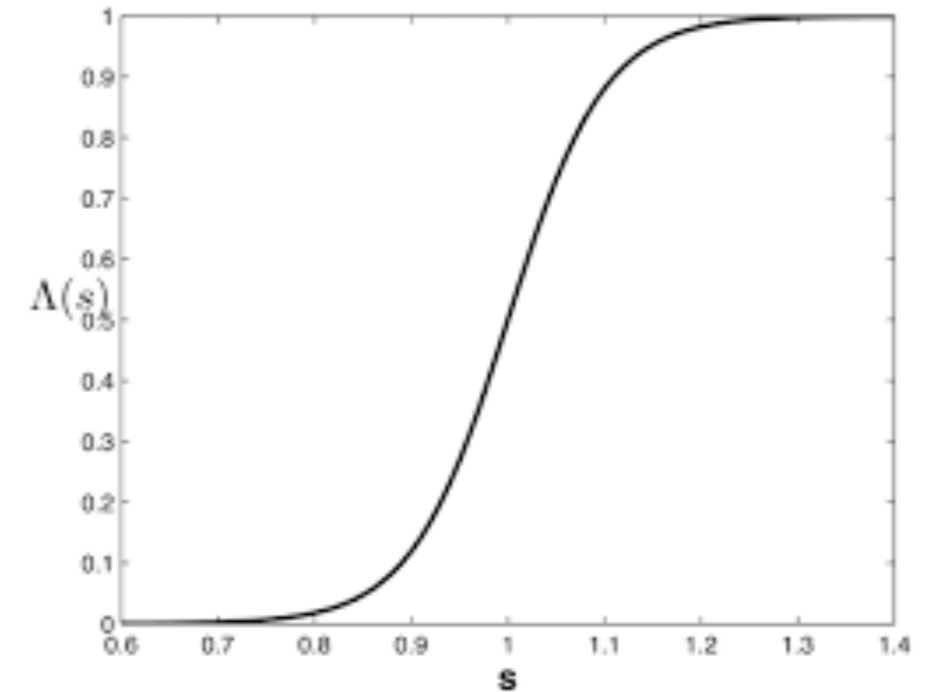
- * i.e. probability of the class times probability of features given the class
- Naïvely assumes features are independent

$$p(c_n | f_1 \dots f_m) = \prod_{i=1}^m p(f_i | c_n) p(c_n)$$

Naïve Bayes

- Pros: Fast to “train” and classify; robust, low-variance; good for low data situations; optimal classifier if independence assumption is correct; extremely simple to implement.
- Cons: Independence assumption rarely holds; low accuracy compared to similar methods in most situations; smoothing required for unseen class/feature combinations

Logistic Regression



- A classifier, despite its name
- A linear model, but uses *softmax* “squashing” to get valid probability

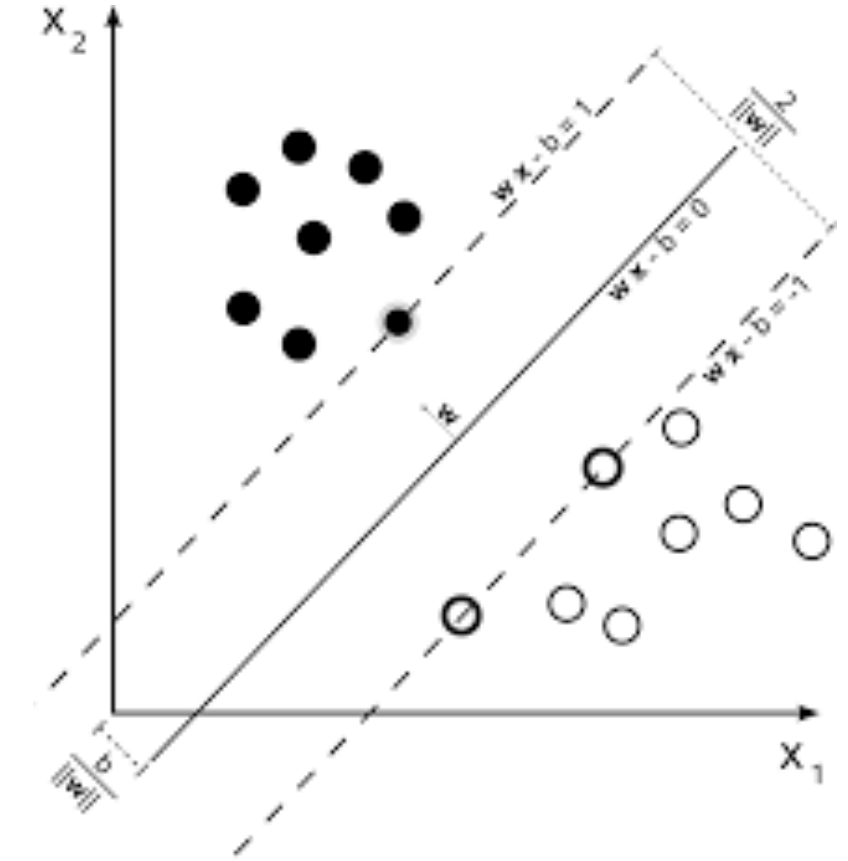
$$p(c_n | f_1 \dots f_m) = \frac{1}{Z} \cdot \exp\left(\sum_{i=0}^m w_i f_i\right)$$

- Training maximizes probability of training data subject to regularization which encourages low or sparse weights

Logistic Regression

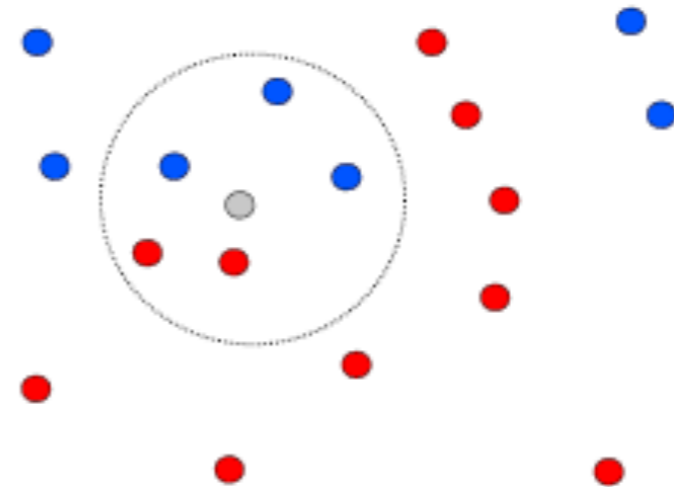
- Pros: A simple yet low-bias classifier; unlike Naïve Bayes not confounded by diverse, correlated features
- Cons: Slow to train; some feature scaling issues; often needs a lot of data to work well; choosing regularisation a nuisance but important since overfitting is a big problem

Support vector machines

- Finds hyperplane which separates the training data with maximum margin
 - * Allows for some misclassification
 - Weight vector is a sum of support vectors (examples on the margin)
- 
- Pros: fast and accurate linear classifier; can do non-linearity with kernel trick; works well with huge feature sets
 - Cons: Multiclass classification awkward; feature scaling can be tricky; deals poorly with class imbalances; uninterpretable

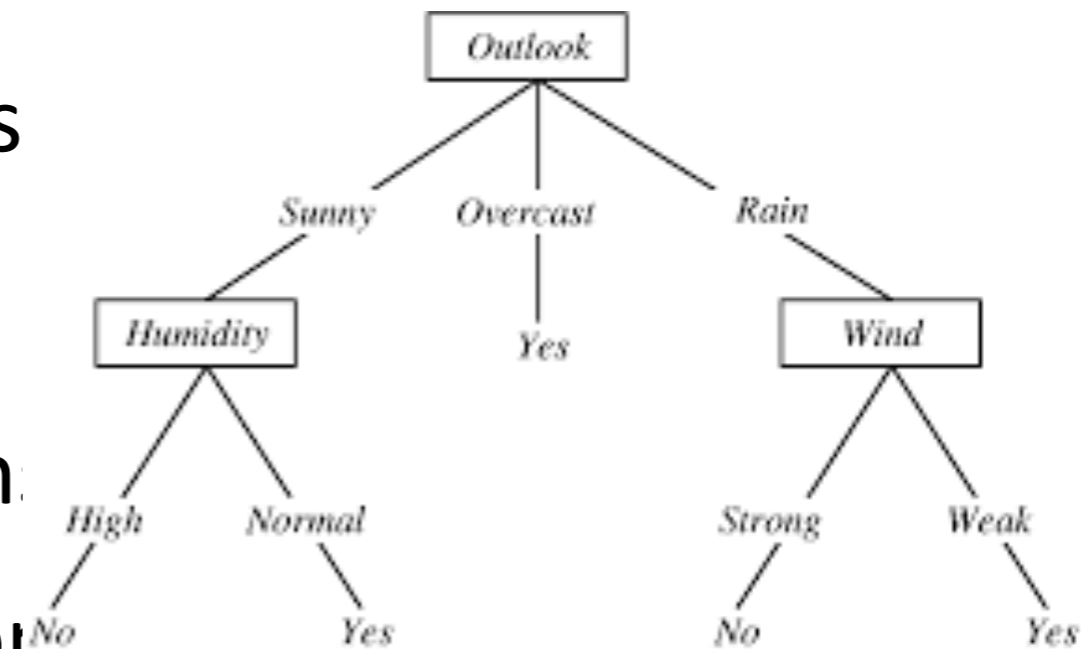
K-Nearest Neighbour

- Classify based on majority class of k -nearest training examples in feature space
- Definition of nearest can vary
 - * Euclidean distance
 - * Cosine distance
- Pros: Simple, effective; no training required; inherently multiclass; optimal with infinite data
- Cons: Have to select k ; issues with unbalanced classes; often slow (need to find those k -neighbours); features must be selected carefully



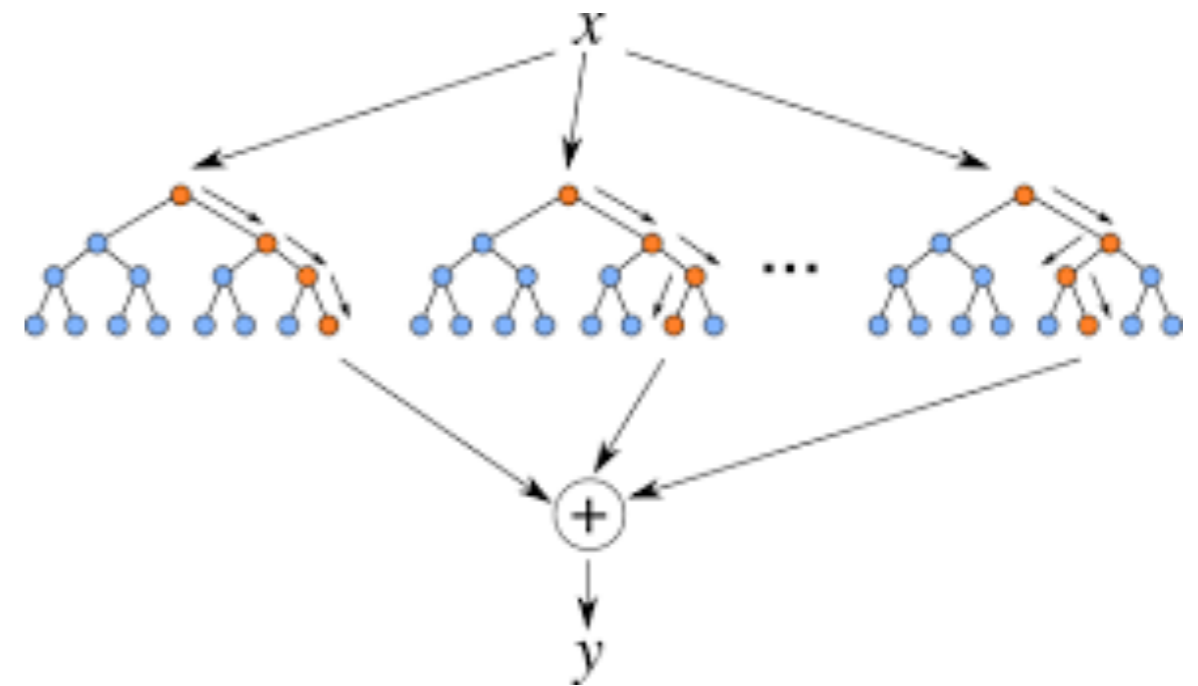
Decision tree

- Construct a tree where nodes correspond to tests on individual features
- Leaves are final class decision
- Based on greedy maximization of mutual information
- Pros: in theory, very interpretable; fast to build and test; feature representation/scaling irrelevant; good for small feature sets, handles non-linearly-separable problems
- Cons: In practice, often not that interpretable; highly redundant sub-trees; not competitive for large feature sets



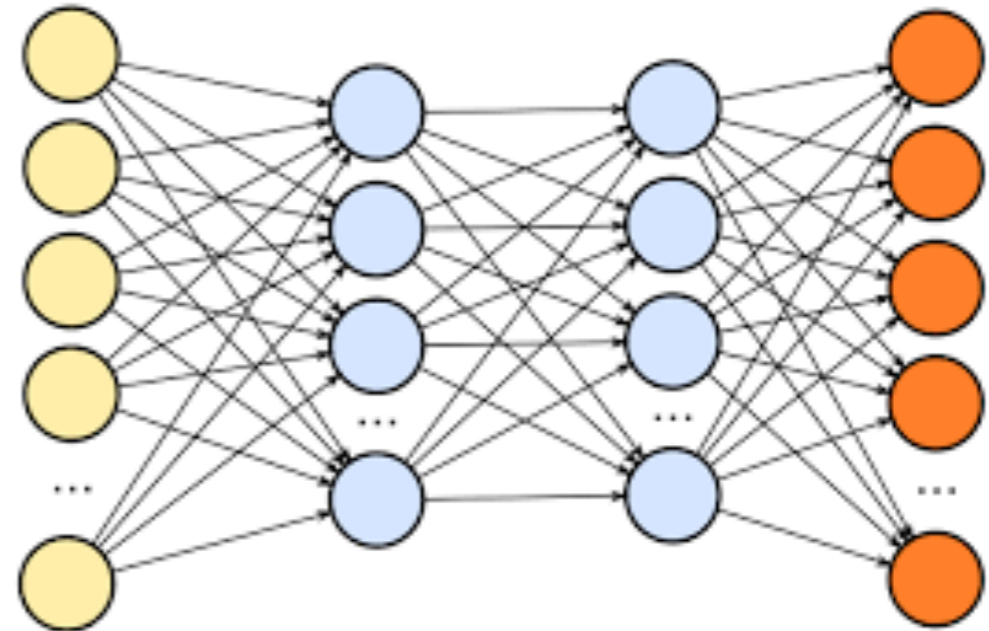
Random forests

- An *ensemble* classifier
- Consists of decision trees trained on different subsets of the training and feature space
- Final class decision is majority vote of sub-classifiers
- Pros: Usually more accurate and more robust than decision trees, a great classifier for small- to moderate-sized feature sets; training easily parallelised
- Cons: Same negatives as decision trees: too slow with large feature sets



Neural Networks

- An interconnected set of nodes typically arranged in layers
- Input layer (features), output layer (class probabilities), and one or more hidden layers



- Each node performs a linear weighting of its inputs from previous layer, passes result through activation function to nodes in next layer
- Pros: Extremely powerful, state-of-the-art accuracy on many tasks in natural language processing and vision
- Cons: Not an off-the-shelf classifier, very difficult to choose good parameters; slow to train; prone to overfitting

HyperParameter tuning

- Dataset for tuning
 - * Development set
 - * Not the training set or the test set
 - * k -fold cross-validation
- Specific hyperparameters are classifier specific
 - * E.g. tree depth for decision trees
- But many hyperparameters relate to regularization
 - * Regularization hyperparameters penalize model complexity
 - * Used to prevent overfitting
- For multiple hyperparameters, use grid search

Evaluation: Accuracy

	Classified As	
Class	A	B
A	79	13
B	8	10

Accuracy = correct classifications/total classifications

$$= (79 + 10)/(79 + 13 + 8 + 10)$$

$$= 0.81$$

0.81 looks good, but most common class baseline accuracy is

$$= (79 + 13)/(79 + 13 + 8 + 10) = 0.84$$

Evaluation: Precision & Recall

Class	Classified As		
	A	B	
A	79	13	False Positives (fp)
B	8	10	True Positives (tp)

False Negatives (fn)

B as “positive class”

Precision = correct classifications of B (tp)
 / total classifications as B (tp + fp)
 = $10 / (10 + 13) = 0.43$

Recall = correct classifications of B (tp)
 / total instances of B (tp + fn)
 = $10 / (10 + 8) = 0.56$

Evaluation: F(1)-score

- Harmonic mean of precision and recall

$$F1 = 2 \text{ precision} * \text{recall} / (\text{precision} + \text{recall})$$

- Like precision and recall, defined relative to a specific positive class
- But can be used as a general multiclass metric
 - * Macroaverage: Average F-score across classes
 - * Microaverage: Calculate F-score using sum of counts

A final word

- Lots of algorithms available to try out on your task of interest (see scikit-learn)
- But if good results on a new task are your goal, then well-annotated, plentiful datasets and appropriate features often more important than the specific algorithm used

Further reading

- J&M3 Ch. 4,5