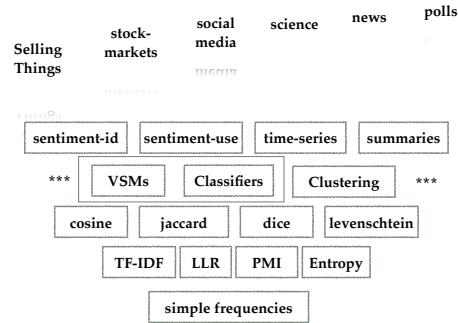


Classification

Lecture 6: Text Analytics for Big Data
Mark Keane, Insight/CSI, UCD



Classification Why We Do This?

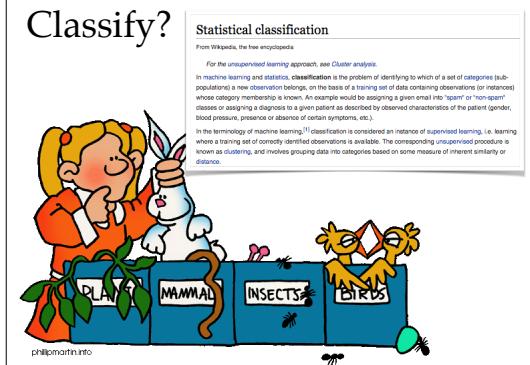
To Classify...

- ◆ Classification is used in text analytics to identify an item as a member of a class and/or to make predictions about that item
- ◆ Is article about sports or business; is this email spam or non-spam; is this buyer profile a binge-shopper
- ◆ If X is a member of class, then I can predict other features; if Mark is a binge-shopper he will need periodic, expanded credit-limits and targeted adverts

Classification is Ubiquitous

- ◆ We have seen similarity, as a way to compare individual items to one another (shopper1 is like shopper 2)
- ◆ Sometimes you have a class/category; so you can say item-X is a category-Y; shopper1 is a BingeShopper or shopper2 is a YoungFamilyShopper
- ◆ Y is a general class of things; BingeShopper is a class
- ◆ Many ML techniques create **classifiers**; programs that can accurately classify items as being in one class or another

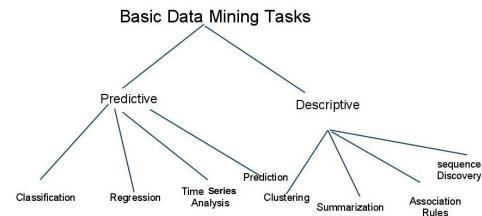
Classify?



Outline

- ◆ Why Do This?
- ◆ Classification (using Supervised Techniques)
- ◆ A Tangent on Cross Validation
- ◆ Next week Clustering...

One View of Tasks...

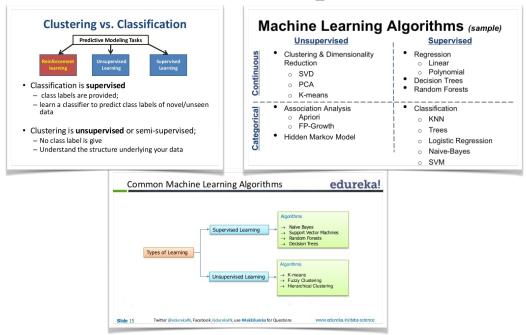


<http://www.slideshare.net/KapilRode/data-mining-32694954>

Classification Examples

- ◆ Based on activity and profiles Tweets have three classes: lurkers, active-tweeters and celebrity-tweeters
- ◆ Company reports can be classified as negative or positive and this predicts share prices
- ◆ Crime-readers buy books by Chandler with "blood" in the title; use this to recommend others
- ◆ Tweeters can be classified as right-wing or left-wing by hashtags, this allows us predict voting preferences

Classification is Supervised...



Selecting Some...

◆ Supervised Learning: Classification

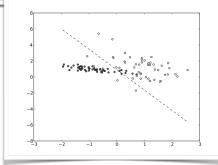
- ◆ k Nearest Neighbours
 - ◆ Naive Bayes
 - ◆ Logistic Regression
 - ◆ Support Vector Machine
- ◆ Unsupervised Learning: Clustering (next lecture)
- ◆ K-Means
 - ◆ Hierarchical Clustering
 - ◆ Graph-based Clustering
- <http://www.slideshare.net/mirjalil/clustering-on-database-systems-rkm-42588852>

Classification Supervised Techniques for Classification

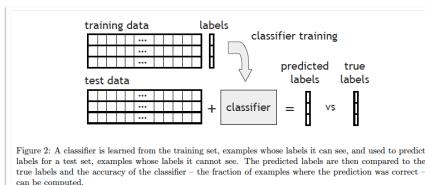
Supervised Techniques

Supervised learning is the machine learning task of inferring a function from labeled training data.^[1] The training data consist of a set of training examples: In supervised learning, each example is a pair consisting of an input object (typically a vector) and a desired output value (also called the supervisory signal). A supervised learning algorithm analyzes the training data and produces an inferred function, which can be used for mapping new examples. An optimal scenario will allow for the algorithm to correctly determine the class labels for unseen instances. This requires the learning algorithm to generalize from the training data to unseen situations in a "reasonable" way (see inductive bias).

Typically, we have a training phase on a sample items, then a test phase in which to test the classifier on unseen items; as in neural networks



Supervised Techniques



Typically, we have a training phase on a sample of items, then a test phase in which to test the classifier on unseen items

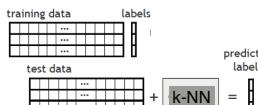
Supervised: Classification

- ◆ k - Nearest Neighbours (k-NN)
- ◆ Naive Bayes
- ◆ Support Vector Machines (SVM)
- ◆ Logistic Regression

Classification Supervised Techniques: k-NN

◆ Training Data: a set of animals each labelled as a panda or a wookie (to find k)

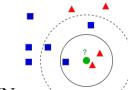
◆ Test Data: a set of unseen animals, with k-NN used to determine if panda/wookie



k -NN: The Idea



- ◆ If you pick 3 most-similar guys to me, and two are pandas and one a wookie, I am probably a panda
 - ◆ If you pick 5 most-similar guys to me and three of them are wookies and two are pandas, then I am probably a wookie
 - ◆ Using k-nearest neighbours; where k is 3 or 5
- <http://en.wikipedia.org/wiki/K-NN>



k -NN: Formulae

Input: D , the set of k training objects, and test object $z = (x', y')$
Process:
 Compute $d(x', x)$, the distance between z and every object, $(x, y) \in D$.
 Select $D_z \subseteq D$, the set of k closest training objects to z .
Output: $y' = \operatorname{argmax}_v \sum_{(x_i, y_i) \in D_z} I(v = y_i)$

Figure 6 provides a high-level summary of the nearest-neighbor classification method. Given a training set D and a test object $z = (x', y')$ the algorithm computes the distance (or similarity) between z and all the training objects $(x, y) \in D$ to determine its nearest-neighbor list, D_z ; (x is the data of a training object, while y is its class. Likewise, x' is the data of the test object and y' is its class.)

Once the nearest-neighbor list is obtained, the test object is classified based on the majority class of its nearest neighbors:

$$\text{Majority Voting: } y' = \operatorname{argmax}_v \sum_{(x_i, y_i) \in D_z} I(v = y_i), \quad (18)$$

where v is a class label, y_i is the class label for the i th nearest neighbors, and $I(\cdot)$ is an indicator function that returns the value 1 if its argument is true and 0 otherwise.

Wu, X. et al. (2008). Top 10 algorithms in data mining. *Knowledge and Information Systems*, 14(1), 1-37.

k -NN: Including Distance



- ◆ Rather than simple voting of neighbours for one class you can measure the distance of items from one another (e.g., cosine similarity / euclidean distance over some tf-idf vector)
- ◆ So, if 10 neighbours have 5 from class-A that are all closer than the 5 from class-B; then weighting the voting on distance will give better classification
- ◆ So, you can make the classification on the basis of the mean distance score for class-A-NNs versus class-B-NNs (nb, this gets more complicated...)

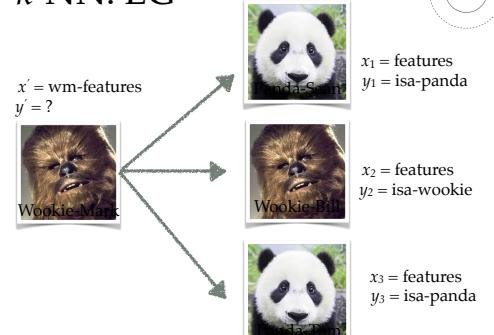
Cunningham, P., & Delany, S. J. (2007). *k*-Nearest neighbour classifiers. *Mult Classif Syst*, 1-17.

k -NN #2: Intruder Detection

Table 1: Analogy between text categorization and intrusion detection when applying the kNN classifier.		
Terms	Text categorization	Intrusion Detection
N	total number of documents	total number of processes
M	total number of distinct words	total number of distinct system calls
n_i	number of times i th word occurs	number of times i th system call was issued
f_{ij}	frequency of i th word in document j	frequency of i th system call in process j
D_j	j th training document	j th training process
X	test document	test process

Liao, Y., & Vemuri, V. R. (2002). Use of K-nearest neighbor classifier for intrusion detection. *Computers & Security*, 21(5), 439-448.

k -NN: EG



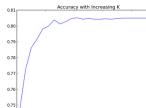
k -NN #1: Text Classification

- ◆ Classic use is in text processing; to take docs and classify them into news categories (sport, current affairs) or wrt medical categories (vascular, psychiatric)
- ◆ Aggressive pre-processing: stop-word removal, stemming, word weighting to achieve an 85% vocabulary reduction
- ◆ Docs then cast as tf-idf vectors, k-NN used
- ◆ Large scale test across several corpora and techniques shows k-NN to do surprisingly well; often used as baseline test

Yang, Y. (1999). An evaluation of statistical approaches to text categorization. *Information retrieval*, 1(1-2), 69-90.

k -NN: The k Issue

Looking at plot, you can see that classifier peaks in accuracy somewhere around 23 neighbors gradually deteriorate. Despite the fact that the classifier maxes out with $k=23$, this doesn't necessarily mean that you should select 23 neighbors.



Take a look at $K=13$. With 13 neighbors, the classifier performs almost just as good as it does with 23 neighbors. In addition to performing nearly as well, it's also interesting to note that K values 15-21 actually have worse performance than $K=13$. This indicates that the classifier is likely to be overfitting, or paying too much attention to the noise in the data. Due to both the comparable performance and the indication of possible overfitting, I would select 13 neighbors for this classifier.

<http://blog.yhatq.com/posts/classification-using-knn-and-python.html>

Main issue in k -NN is knowing k ; different ks may give different answers, you want most accurate; often done empirically

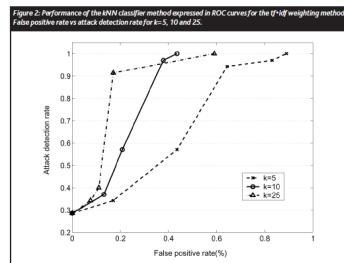
k -NN #2: Intruder Detection

Process ID 994									
close	error	open	mmap	open	mmap	mmmap	mmmap	mmmap	mmmap
mmap	close	open	mmap	close	open	mmap	mmap	mmap	mmap
mmap	close	close	mmmap	open	iced	access	closed	iced	iced

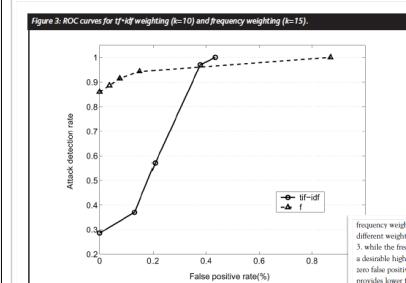
- ◆ 1998 DARPA BSM Competition; idea is to capture "normal" behaviour of users, systems, networks, match new behaviour to these and classify as normal/abnormal
- ◆ Uses program profiles broken down in the freqs of system calls (unix commands: *su chmod open close*)
- ◆ Corpus is a set of programs, which contain some known "attacks" on various US Navy Systems

Liao, Y., & Vemuri, V. R. (2002). Use of K-nearest neighbor classifier for intrusion detection. *Computers & Security*, 21(5), 439-448.

k -NN #2: Intruder Detection



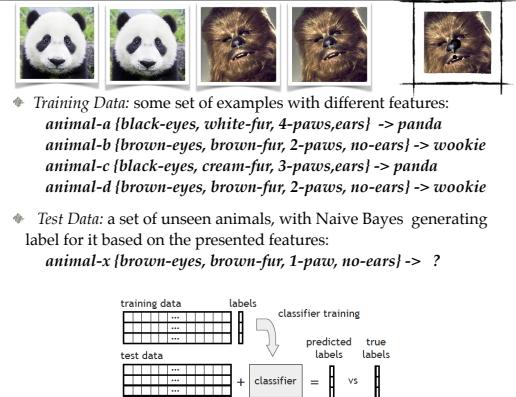
Liao, Y., & Vemuri, V. R. (2002). Use of K-nearest neighbor classifier for intrusion detection. *Computers & Security*, 21(5), 439-448.



Liao, Y., & Vemuri, V. R. (2002). Use of K-nearest neighbor classifier for intrusion detection. *Computers & Security*, 21(5), 439-448.

frequency weight. A comparison of two different weighting methods is shown in Figure 3 while the frequency weighting method offers a desirable high attack detection rate (80%) at zero false positive rate. The $t^{\alpha}df$ weighting method provides lower false positive rate at 10% attack detection rate. It appears that the $t^{\alpha}df$ weighting can make process vectors of two classes dissimilar even with the same frequency weighting. Therefore, a lower threshold value is needed, and better false positive rate can be achieved with the $t^{\alpha}df$ weighting method.

Classification
Supervised Techniques:
Naive Bayes



Naive Bayes:
The Idea



◆ "Most" of the wookies I have seen have brown eyes and pandas have black; one or two wookies have black eyes, and I saw one panda with brown eyes (but pandas can also have blue, green and red eyes!).

◆ Use the probability of features (e.g., eye-colour) occurring in some sample of instances (wookies and pandas I have seen), wrt a category (wookie or panda), to predict the probability than an unseen Beast-X is in one category or the other

◆ Clever Bayes bit tells you how to combine the probabilities of different features, to compute an overall probability for X

http://en.wikipedia.org/wiki/Naive_Bayes_classifier

Naive Bayes: Formula

Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $\mathbf{x} = (x_1, \dots, x_n)$ representing some n features (dependent variables), it assigns to this instance probabilities $p(C_k | x_1, \dots, x_n)$ for each of k possible outcomes or classes.^[7]

The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number of values, then basing such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable. Using Bayes' theorem, the conditional probability can be decomposed as

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

In plain English, using Bayesian probability terminology, the above equation can be written as
 posterior = prior \times likelihood / evidence.

http://en.wikipedia.org/wiki/Naive_Bayes_classifier

Naive Bayes: Bayes Theorem

- Bayesian classifiers use **Bayes theorem**, which says

$$p(c_j | d) = \frac{p(d | c_j) p(c_j)}{p(d)}$$

- $p(c_j | d)$ = probability of instance d being in class c_j
 This is what we are trying to compute
- $p(d | c_j)$ = probability of generating instance d given class c_j
 We can imagine that being in class c_j causes you to have feature d with some probability
- $p(c_j)$ = probability of occurrence of class c_j
 This is just how frequent the class c_j is in our database
- $p(d)$ = probability of instance d occurring

This can actually be ignored, since it is the same for all

http://www.cs.ucr.edu/~eamonn/CE/Bayesian%20Classification%20withInsect_examples.pdf

Naive Bayes: EG (single feature)

Eyes	Beast
brown	panda
black	wookie
brown	wookie
brown	wookie
red	panda
brown	wookie
blue	wookie
green	panda

...

Naive Bayes: EG (single feature)

Eyes	Beast
brown	panda
black	wookie
brown	wookie
brown	wookie
red	panda
brown	wookie
blue	wookie
green	panda

$p(c_j | d) = \frac{p(d | c_j) p(c_j)}{p(d)}$

$p(\text{brown/panda}) = 1/3$
 $p(\text{panda}) = 3/8$
 $p(\text{brown}) = 4/8$

$p(\text{panda/brown}) = \frac{1/3 * 3/8}{4/8} = 0.125$
 $4/8$

$p(\text{wookie/brown}) = \frac{3/5 * 5/8}{4/8} = 0.375$
 $4/8$

Naive Bayes: EG (many features)

Officer Drew and Nina McGrew....

Name	Over 170cm	Eye	Hair length	Sex
Drew	No	Blue	Short	Male
Claudia	Yes	Brown	Long	Female
Drew	No	Blue	Long	Female
Alberto	Yes	Brown	Short	Male
Karin	No	Blue	Long	Female
Nina	Yes	Brown	Short	Female
Sergio	Yes	Blue	Long	Male

http://www.cs.ucr.edu/~eamonn/CE/Bayesian%20Classification%20withInsect_examples.pdf

Naive Bayes: Turtles !

- To simplify the task, **naive Bayesian classifiers** assume attributes have independent distributions, and thereby estimate

$$p(d | c_j) = p(d_1 | c_j) * p(d_2 | c_j) * \dots * p(d_n | c_j)$$

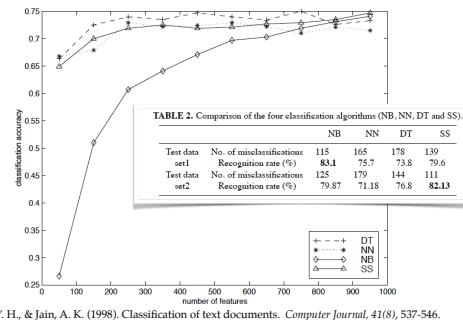
The probability of class c_j generating instance d , equals...

The probability of class c_j generating the observed value for feature 1, multiplied by...

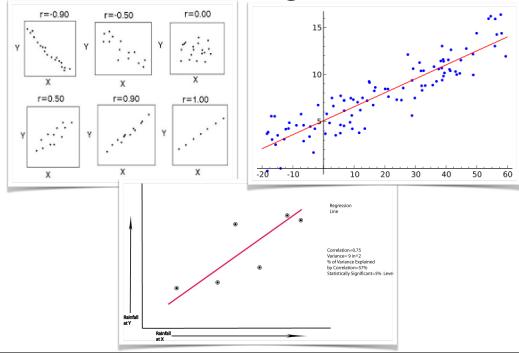
The probability of class c_j generating the observed value for feature 2, multiplied by...

http://www.cs.ucr.edu/~eamonn/CE/Bayesian%20Classification%20withInsect_examples.pdf

Naive Bayes: Li & Jain (1998)



Correlations & Regression: Pics



Logistic Regression: Formula

$$\ln\left(\frac{P}{1-P}\right) = a + bX$$

formula relates the log of the odds ($P/1-P$) to the linear function for the explanatory variable

$$\frac{P}{1-P} = e^{a+bX}$$

a is intercept, b the slope of the regression line and X the value of the independent variable

$$P = \frac{e^{a+bX}}{1+e^{a+bX}}$$

when there are more variables you just add more terms to the linear fn

[http://en.wikipedia.org/wiki/Jaccard_index](https://en.wikipedia.org/wiki/Jaccard_index)

Logistic Regression: Explanation

Definition of the logistic function [edit]
 Logistic regression is based on an application of the logistic function. The logistic function is used because it can take an input with any value from negative to positive infinity, whereas the output always takes values between zero and one^[1] and hence is interpretable as a probability. The logistic function is defined as follows:

$$\sigma(t) = \frac{e^t}{1+e^t} = \frac{1}{1+e^{-t}}$$

A graph of the logistic function is shown in Figure 1.

If t is replaced as a linear function of an explanatory variable x (or of a linear combination of explanatory variables), then we get the form of the logistic function:

$$t = \beta_0 + \beta_1 x$$

And the logistic function can now be written as:

$$\sigma(x) = \frac{e^{x-\beta_0-\beta_1 x}}{1+e^{x-\beta_0-\beta_1 x}}$$

Note that $\sigma(x)$ is interpreted as the probability of the dependent variable equaling a "success" or "lose" rather than a failure or non-case. It's clear that the response variables y_i are not identically distributed: $P(Y_i = 1 | X)$ differs from one data point, X_i , to another, though they are independent given design matrix X and shared with parameters β .

Definition of the inverse of the logistic function [edit]

We can now define the inverse of the logistic function, y , the logit (log odds):

$$g(F(x)) = \ln \frac{F(x)}{1-F(x)} = \beta_0 + \beta_1 x,$$

and equivalently

$$\frac{F(x)}{1-F(x)} = e^{\beta_0+\beta_1 x}$$

Interpretation of these terms [edit]

In the above equation, the terms are as follows:

- β_0 is referred to as the logit function. The graph of $g(F(x))$ illustrates that the logit (i.e., log odds or natural logarithm of the odds) is equivalent to the linear regression expression.
- β_1 is referred to as the residual logit.
- $\beta_0 + \beta_1 x$ is the term that the dependent variable equals a case, given some linear combination of the predictors. The formula for $F(x)$ illustrates that the probability of the dependent variable equaling a case is equal to the value of the logistic function of the linear regression expression. This is important in that it shows that the value of the linear regression expression can vary from negative to positive infinity and yet, after transformation, the resulting expression for the logistic $F(x)$ will be bounded between 0 and 1.
- β_1 is the regression coefficient multiplied by some value of the predictor.
- β_0 is the natural logarithm.

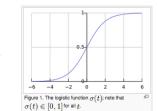
https://en.wikipedia.org/wiki/Logistic_regression

Classification Supervised Techniques: Logistic Regression

Correlations & Regression

- ◆ **Correlation:** tells you how two sets of figures (two variables) co-vary with one another; using co-variance and variance (i.e., degree of linear dependence between 2 vars)
- ◆ **Linear Regression:** tells you about the relationship between some dependent variable and one or more explanatory variables (>1 being Multiple Regression)

https://en.wikipedia.org/wiki/Linear_regression



Regression: The Idea

- ◆ **Logistic Regression** projects "normal" regression into a logistic function, so that the output is between 0 and 1 (not cont.)
- ◆ It gives the probability that something will/will not (1/0) occur (or be the case) given value(s) of some independent variable(s)

https://en.wikipedia.org/wiki/Logistic_regression

Logistic Regression: EG1

- ◆ **Data:** Imagine you have some data on the number of servers that are hacked (1) or not hacked (0) after different time periods in weeks (28 to 33 weeks)
- ◆ **Logistic Regression:** will give you a probability for whether a new server is likely to be hacked, based on this data, after a given time period (say 31 weeks)

<http://vassarstats.net/logreg1.html>

Logistic Regression: EG1

Data: The number of servers that are hacked (1) or not hacked (0) after different time periods measured in weeks (28 to 33 weeks)

I	II	III	IV	V	VI	VII
X	0	1	II+III	Total	Y as Observed Probability	Y as Odds
28	4	2	6	.3333	.5000	.6931
29	3	2	5	.4000	.6667	.4055
30	2	7	9	.7778	.3000	1.2528
31	2	7	9	.7778	3.0000	1.2528
32	4	16	20	.8000	4.0000	1.3863
33	1	14	15	.9333	14.0000	2.6391

$$\text{Odds} = \frac{\text{probability of the event}}{1 - \text{probability of event}} = \frac{P}{1-P}$$

OR

$$P = \frac{\text{Odds}}{1 + \text{Odds}}$$

<http://vassarstats.net/logreg1.html>

For each level of X, the weighting factor is the number of observations for that level.

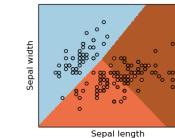
Intercept=-17.2086 is the point on the Y-axis (log odds) crossed by the regression line when X=0.
Slope=.5934 is the rate at which the predicted log odds increases (or, in some cases, decreases) with each successive unit of X. Within the context of logistic regression, you will usually find the slope of the log odds regression line referred to as the "constant."

The exponent of the slope: $\exp(.5934) = 1.81$ describes the proportionate rate at which the predicted odds changes with each successive unit of X. In the present example, the predicted odds for X=29 is 1.81 times as large as the one for X=28; the one for X=30 is 1.81 times as large as the one for X=29; and so on.

<http://vassarstats.net/logreg1.html>

Logistic Regression 3-class Classifier

Show below is a logistic-regression classifiers decision boundaries on the iris dataset. The datapoints are colored according to their labels.



Logistic Regression: Prog

```
Python source code: plot_iris_logistic.py
print(__doc__)

# Code source: Gaël Varoquaux
# Modified for documentation by Jaques Grobler
# License: BSD 3 clause

import numpy as np
import matplotlib.pyplot as plt
from sklearn import linear_model, datasets

# import some data to play with
np.random.seed(0)
X = iris.data[:, :2] # we only take the first two features.
y = iris.target
h = .02 # step size in the mesh

logreg = linear_model.LogisticRegression(C=1e5)
logreg.fit(X, y)

# Plot the decision boundary. For that, we will assign a color to each
# point in the mesh [x_min, x_max]x[y_min, y_max].
x_min, x_max = X[:, 0].min(), X[:, 0].max()
y_min, y_max = X[:, 1].min(), X[:, 1].max()
xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
                     np.arange(y_min, y_max, h))

Z = logreg.predict(np.c_[xx.ravel(), yy.ravel()])

# Put the result into a color plot
Z = Z.reshape(xx.shape)
plt.figure(1, figsize=(4, 3))
plt.pcolormesh(xx, yy, Z, cmap=plt.cm.Paired)

# Plot also the training points
plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Paired,
            edgecolors='black', s=15)
plt.title('Logistic Regression 3-class Classifier')
plt.xlabel('Sepal length')
plt.ylabel('Sepal width')

plt.show()

http://scikit-learn.org/stable/auto_examples/linear_model/plot_iris_logistic.html
```

LR: Genres



- ◆ Genre: “reflects way text created, how it was distributed, register of language, audience addressed”
- ◆ So, may be about words used, expressions, syntactic constructions, ...
- ◆ **Generic Facets:** features of text identified as important in determining genre: *broadcast* V *directed*, *narrative* / *susasive* / *descriptive* / *directive*

Kessler, B., Numberg, G., & Schütze, H. (1997, July). Automatic detection of text genre. In Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (pp. 32-38). ACL.

LR: Genre Facets

- ◆ Facets are quite high-level based on lower-level “visible” features:
 - *structural*: passives, nominalisation, POS counts
 - *lexical*: terms of address Mr...Mrs (register)
 - *character-level*: punctuation, !s, ?s
 - *derivative*: ratios of words, sentence lengths...

Kessler, B., Numberg, G., & Schütze, H. (1997, July). Automatic detection of text genre. In Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (pp. 32-38). ACL.

LR: Genre Facets

- ◆ 499 texts Brown Corpus: training (402), test (97)
- ◆ Facets: Brow (popular, middle...high), Narrative (yes or no), Genre (reportage, editorial, scitech, legal, fiction, nonfiction)
- ◆ Use LR for each facet on each training text; then tested against unseen 97 (NN compare)

Kessler, B., Numberg, G., & Schütze, H. (1997, July). Automatic detection of text genre. In Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (pp. 32-38). ACL.

LR: Results

Genre	Baseline		LR (Surf.)		2LP		3LP		LR (Struct.)	
	All	Sel.	All	All	All	All	All	Sel.	All	Sel.
Rep	81	89*	88	94*	94*	90*	90*	90*	90*	90*
Edit	81	75	75	74	80	79	77			
Legal	95	96	96	95	95	93	93			
Scitech	94	100*	96	99*	94	93	96			
Nonfict	67	67	68	78*	67	73	74			
Fict	81	93*	96*	99*	81	96*	96*			
Brow										
Popular	74	74	75	74	74	72	73			
Middle	68	66	67	64	64	58	64			
Uppermiddle	88	74	78	86	88	79	82			
High	70	84*	88*	89*	90*	85*	86*			

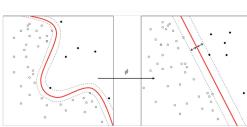
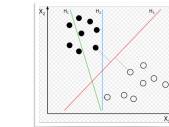
Note: Numbers are the percent of the evaluation subcorpus ($N = 97$) which was correctly classified on a binary discrimination task. The Baseline column tells what percentage would be got correct by guessing No for each level. Headers have the same meaning as in Table 1.
* means significantly better than Baseline at $p < .05$, using a binomial distribution ($N=97$, p as per first column).

Kessler, B., Numberg, G., & Schütze, H. (1997, July). Automatic detection of text genre. In Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (pp. 32-38). ACL.

Classification Supervised Techniques: Support Vector Machines

SVM: The Idea

- ◆ Explores hyperplanes for an multi-dimensional space; finding the “best” one; geometrically, the one with the largest separation (or margin) between classes
- ◆ Basically, linear classification but can use **kernel** idea to handle non-linear classifications (projecting into linear)



In geometry a hyperplane is a subspace of one dimension less than its ambient space. If a space is 3-dimensional then its hyperplanes are the 2-dimensional planes, while if the space is 2-dimensional, its hyperplanes are the 1-dimensional lines. This notion can be used in any general space in which the concept of dimension of a subspace is defined.

SVM: The Idea

H1 does not separate classes
H2 does, but with small margin
H3 separates with maximal margin

http://en.wikipedia.org/wiki/Support_vector_machine

SVM: Formula

The reason why SVM insists on finding the maximum margin hyperplanes is that it offers the best generalization ability. It allows not only the best classification performance (e.g., accuracy) on the training data, but also leaves much room for the correct classification of the future data. To ensure that the maximum margin hyperplanes are actually found, an SVM classifier attempts to maximize the following function with respect to $\hat{\mathbf{w}}$ and b :

$$L_P = \frac{1}{2} \|\hat{\mathbf{w}}\|^2 - \sum_{i=1}^t \alpha_i y_i (\hat{\mathbf{w}} \cdot \mathbf{x}_i + b) + \sum_{i=1}^t \alpha_i \quad (2)$$

where t is the number of training examples, and $\alpha_i, i = 1, \dots, t$, are non-negative numbers such that the derivatives of L_P with respect to α_i are zero. α_i are the Lagrange multipliers and L_P is called the Lagrangian. In this equation, the vectors $\hat{\mathbf{w}}$ and constant b define the hyperplane.

SVM: Formula

Linear SVM (soft)
Given point learning data D a set of n points of the form $\mathcal{P} = \{(\mathbf{x}_i, y_i) | \mathbf{x}_i \in \mathbb{R}^d, y_i \in \{-1, 1\}\}_{i=1}^n$ where the y_i is either $+1$ or -1 , indicating the class to which the point \mathbf{x}_i belongs. Each \mathbf{x}_i is a d -dimensional real vector. We want to find the maximum margin hyperplane that divides the points having $y_i = +1$ from those having $y_i = -1$. Any hyperplane can be written as the set of points \mathbf{x} satisfying

$$\mathbf{w} \cdot \mathbf{x} + b = 0,$$

where \mathbf{w} denotes the dot product and \mathbf{w} the (not necessarily normalized) normal vector to the hyperplane. The parameter $\|\mathbf{w}\|$ determines the offset of the hyperplane from the origin along the normal vector \mathbf{w} .

If the training data are linearly separable, we can select two hyperplanes in a way that they separate the data and there are no points between them, and then try to maximize their distance. The region bounded by them is called ‘the margin’. These hyperplanes can be defined by the equations

$$\mathbf{w} \cdot \mathbf{x}_i + b = 1 \quad \text{for } \mathbf{x}_i \text{ of the first class}$$

and

$$\mathbf{w} \cdot \mathbf{x}_i + b = -1. \quad \text{for } \mathbf{x}_i \text{ of the second class}$$

By using geometry, we find the distance between these two hyperplanes is $\frac{2}{\|\mathbf{w}\|}$, so we want to minimize $\|\mathbf{w}\|$. As we also have to prevent data points from falling into the margin, we get the following constraint: for each i either

$$\mathbf{w} \cdot \mathbf{x}_i + b \geq 1 \quad \text{for } \mathbf{x}_i \text{ of the first class}$$

or

$$\mathbf{w} \cdot \mathbf{x}_i + b \leq -1 \quad \text{for } \mathbf{x}_i \text{ of the second class}$$

This can be rewritten as:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1, \quad \text{for all } 1 \leq i \leq n. \quad (1)$$

We can put this together to get the optimization problem:

$$\begin{aligned} & \text{Minimize}_{(\mathbf{w}, b)} && \|\mathbf{w}\| \\ & \text{subject to for all } i = 1, \dots, n, && y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1. \end{aligned}$$

http://en.wikipedia.org/wiki/Lagrangian_index

SVM: More...

- ◆ Linear SVM assumes linear separable
- ◆ But, there are extensions using kernel function to non-linear cases; many different including quadratic and exponential kernel functions

In Euclidean geometry linear separability is a geometric property of a pair of sets of points. This is most easily visualized in two dimensions (the Euclidean plane) by thinking of one set of points as being colored blue and the other set of points as being colored red. These two sets are linearly separable if there exists at least one line in the plane with all the blue points on one side of the line and all the red points on the other side. This idea immediately generalizes to higher-dimensional Euclidean spaces if line is replaced by hyperplane.

Whereas the original problem may be stated in a finite dimensional space, it often happens that the sets to discriminate are not linearly separable in that space. For this reason, it was proposed that the original finite-dimensional space be mapped into a much higher-dimensional space, presumably making the separation easier in that space. To keep the computational load reasonable, the mappings used by SVM schemes are designed to ensure that dot products may be computed easily in terms of the variables in the original space, by defining them in terms of a kernel function $k(x, y)$ selected to suit the problem.^[2]

SVM Eg: SPAM:Context

- ◆ 2014, 4.1B email accounts; 196.3B e-mails per day (108.7B business); 12% of business mails are spam (after spam filtering); 13B spam mails... (www.radicati.com “email statistics”)
- ◆ Traditional methods: white-lists, black-lists, handcrafted rules for text (“viagra”), patterns like missing headers, body only pic, too many non-words (v.i.a.g.r.a)
- ◆ Idea: to have tractable learning mechanism for spam/non-spam
- ◆ Nearly every classifier has been applied to spam

Androutsopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. “DEMOKRITOS”, National Center for Scientific Research.

SVM Eg: SPAM:The Problem

- ◆ Take a corpus of emails, labelled as spam/non-spam; create a classifier to distinguish two classes, apply to incoming mails
- ◆ Spam-detection becomes a text-classification problem; tailored to specific mails of a given users
- ◆ Also, a cost-sensitive problem; blocking a non-spam mail has higher cost than allowing a spam mail (weighted classification)
- ◆ Also, look at encoding issues; size of training set, size of attribute set, use of ngrams as features and Sequential Minimal Optimisation (SMO), time and memory efficient version of SVM

Androutsopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. “DEMOKRITOS”, National Center for Scientific Research.

SVM Eg: SPAM: Data

- ◆ Training works on smallish sets; 100s and 1000s of emails; concentrated on non-“regular correspondents”
- ◆ Headers, html removed, but stop-words left, no lemmatisation and “\$” and “£” treated as tokens (nb stops in twitter)
- ◆ Each message is a feature-vector; sometimes binary, but they used frequency and length of doc (“\$”, “money”)
- ◆ Also ngrams as they capture typical phrases: “get rich”, “rich quick”; tho’ with size of n, this grows exponentially, no $n > 3$

Androutsopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. “DEMOKRITOS”, National Center for Scientific Research.

SVM Eg: SPAM: Subplot

- ◆ ngrams will produce many features and sparsity; so how do you cut them down
- ◆ Took all items; 1gram, 2grams, 3grams
- ◆ Remove all items with count < 4;
- ◆ Get Information Gain of every feature wrt its category (legitimate v spam); reduces entropy/redundancy of the feature set; optimises difference in features wrt category

Androutsopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. “DEMOKRITOS”, National Center for Scientific Research.

SVM Eg: SPAM: Feature Selection

Table II. The n-grams with the highest information gain scores in PUI, for $n \in \{1, 2, 3\}$, ranked by decreasing information gain score. The last and conditional probabilities of the n-grams, assuming Boolean attribute. Shaded and non-shaded rows indicate n-grams that are more frequent in the spam and legitimate category, respectively.

n	n-gram	$P(X_1 = 1)$	$P(X_1 = 1 c_S)$	$P(X_1 = 1 c_G)$
1	fax :	0.174386	0.090677	0.020270
1	university	0.089800	0.747193	0.022774
1	click	0.118960	0.922131	0.022131
1	mailing	0.169845	0.29932	0.351966
1	natural	0.156221	0.277719	0.020270
1	the internet	0.089806	0.010112	0.252587
1	newspaper	0.115988	0.012232	0.252587
1	research	0.281562	0.451561	0.072463
1	university of	0.154405	0.274193	0.020270

Androutsopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. “DEMOKRITOS”, National Center for Scientific Research.

SVM Eg: SPAM: Complexity

Table III. Computational complexity of the four learning algorithms. N is the number of training messages, $m \leq M$ the number of selected attributes, and m' the number of iterations in LogitBoost.

algorithm	training	classification
Naive Bayes	$O(mN)$	$O(m)$
Flexible Bayes	$O(mN)$	$O(mN)$
SVM with SMO	$O(mN^2)$	$O(mN)$
linear SVM with SMO	$O(mN^2)$	$O(m)$
LogitBoost with regression stumps	$O(m'mN^2)$	$O(m')$

Androustopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. " DEMOKRITOS", National Center for Scientific Research.

SVM Eg: SPAM: Results

Table VIII. Real-life evaluation results of Filtron, using the SVM with 520 1-gram attributes, for $\lambda = 1$. The SVM was trained on PC3. Bracketed precision, recall, and WAcc scores are the corresponding scores we had obtained with 10-fold cross-validation on PC3.

	messages received	(avg. 31.75 per day)
spam messages received	6732	(avg. 7.66 per day)
legitimate messages received	5109	(avg. 24.10 per day)
legitimate-to-spam ratio	3.15	
correctly classified legitimate messages ($L \rightarrow L$)	5027	
incorrectly classified legitimate messages ($L \rightarrow S$)	52	(avg. 1.72 per week)
correctly classified spam messages ($S \rightarrow S$)	1450	
incorrectly classified spam messages	173	(avg. 5.71 per week)
precision	96.54% (n=3: 96.43%)	
recall	89.34% (n=3: 95.05%)	
WAcc	96.66% (n=3: 96.22%)	

Of the 52 fails, 52% were auto-mails (newsletters), 22% 2-3 word mails with attachments, 26% were causal 1-2 line mails in spam style with no attachments/hyperlinks (lambda relates to cost fn)

Androustopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. " DEMOKRITOS", National Center for Scientific Research.

SVM: Issues

- Theoretically well-founded and can work with a few dozen examples, also insensitive to number of dimensions; but nb complexity



Classification A Tangent on Cross-Validation

A Problem ?

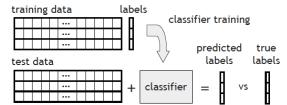


Figure 2: A classifier is learned from the training set, examples whose labels it can see, and used to predict labels for a test set, examples whose labels it cannot see. The predicted labels are then compared to the true labels and the accuracy of the classifier – the fraction of examples where the prediction was correct – can be computed.

- Note, we have built this classifier on a sample of training data (randomly selected?)
- It might not be a good/representative sample making classifier crap/overfitted

A Solution?

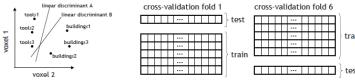


Figure 3: Left: Learning a linear classifier is equivalent to learning a line that separates examples in the two classes (vectors in a 2-D visual) as well as possible. Right: During cross-validation each of 6 groups of examples takes a turn as the test set while the rest serve as the training set.

- Cross-validation deals with this problem
- Use several different (systematic) samples of training sets and report the mean solution
- Also indicates the model's performance; very diff answers for diff training sets is bad

Types of Cross Validation

- Exhaustive:** (nb complexity)
 - leave-p-out: use p items as test subset; the rest as training subset; explore all subsets of p-items
 - leave-one-out: same, but where p = 1
- Non-Exhaustive:** (nb better on complexity)
 - k-fold: randomly partition items into k subsamples; and use one for test, rest for training
 - 2-fold: same where k is 2
 - repeated random sub-sampling: as it suggests

k-fold cross-validation

- Partition original sample into k sub-samples, by random selection of items
- One sub-sample is used as test-subset, and remaining $k-1$ sub-samples as training subset
- Cross-validation is repeated k times for each; with results combined or averaged
- Stratified, k-fold; mean response value is kept approx. equal; e.g. 50/50 split in binary label set
- Benefit: all training set is used and better complexity than exhaustive methods

k-fold cross-validation: EG

Androustopoulos, I., Palioras, G., & Michelakis, E. (2004). Learning to filter unsolicited commercial e-mail. " DEMOKRITOS", National Center for Scientific Research.

All of the experiments were performed using stratified ten-fold cross-validation. That is, each message collection was divided into ten equally large parts, maintaining the original distribution of the two classes in each part. Each experiment was repeated ten times, reserving a different part for testing at each iteration, and using the remaining parts for training. The results were then averaged over the ten iterations, producing more reliable results, and allowing the entire corpus to be exploited for both training and testing. The attribute selection process was repeated anew at each iteration of the cross-validation. When we refer to the training corpus of cross-validation experiments, we mean the nine parts that were used for training at each iteration.

2-fold cross-validation

1. Partition original sample into 2 equally-sized sub-samples, by random selection of items (usually shuffle the data a split it in two)
2. So, you have two subsets; one the test-set, the other the training subset
3. Cross-validation is repeated *twice*, swapping these subsets as test/training
 - Benefit: all items used for either training or testing; subsets are large

Repeated, random sub-sampling

- ◆ Randomly split the sample into training and testing subsets
- ◆ Repeat this a number of times; taking average of results
- ◆ NB: may mean a given item never appears in training set

Classifications
Conclusions

Ahem...

- ◆ Classification and the use of classifiers is a standard solution to lots of problems
- ◆ Given some data-set; can I just develop one of these programs and then re-use it
- ◆ There are, however, many issues about their use that need to be kept in mind

Issues...

- ◆ Classifiers are only as good as your data-set; if it is unrepresentative or biased then it may be learning weird stuff (tanks+skies)
- ◆ When something is learned, we may not really know what it has learned (e.g., what are the important features; could be critical)
- ◆ There are dark arts and key choices required to make these work (e.g., pre-processing)