Colin Carroll

An Overview of Classification Algorithms.

Colin Carroll

June 12, 2013

Colin Carroll An Overview of Classification Algorithms.

Colin Carroll

Goals

- Overview of what the problem is.
- High level view of what is available.
- Some intuition of how it works.

Colin Carroll

Goals

3 2 ч 3 0 8 e 8 0 ລ 06 2 Э 3 S 0 3 9 0 2 208

Colin Carroll

Not goals

- Feature detection.
- Dimensionality reduction.
- Tuning models.
- Pros and cons.



Colin Carroll



e python



Colin Carroll

The iris dataset

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
51	7.0	3.2	4.7	1.4	versicolor
52	6.4	3.2	4.5	1.5	versicolor
101	6.3	3.3	6.0	2.5	virginica
102	5.8	2.7	5.1	1.9	virginica



Colin Carroll



	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
51	7.0	3.2	4.7	1.4	versicolor
52	6.4	3.2	4.5	1.5	versicolor
101	6.3	3.3	6.0	2.5	virginica
102	5.8	2.7	5.1	1.9	virginica

Colin Carroll

The iris dataset



Colin Carroll

Supervised vs. unsupervised learning





Colin Carroll

Unsupervised learning

Should expect results to be worse.

- Not magic.
- Easiest if we assume all data is 2 dimensional.

Colin Carroll

K-means clustering

We assume there will be k clusters. Initialize by selecting k points μ_1, \ldots, μ_k . Proceed by

Assigning each point to the μ_j nearest to it. Gives k sets C_1, \ldots, C_k .

■ Choose new μ_j 's by letting μ_j = average(C_j). Stop when the C_i 's stop changing.

Colin Carroll

K-means clustering

Given N observations $X = \{x_1, \ldots, x_N\}$, find a partition S_1, \ldots, S_K so that

$$S_1,\ldots,S_K = \arg\min\sum_{i=1}^K \sum_{x_j \in S_i} \|x_j - \mu_i\|.$$

Colin Carroll



Colin Carroll

	1	2	3
setosa	50	0	0
versicolor	0	12	38
virginica	0	35	15

Colin Carroll



Colin Carroll

	1	2	3
setosa	50	0	0
versicolor	0	2	48
virginica	0	36	14

Colin Carroll

Mixture Modeling

We assume the members of each (unknown) cluster has been drawn from a different population.



Colin Carroll

Mixture Modeling



Colin Carroll

Mixture Modeling



Colin Carroll

Mixture Modeling

	1	2	3
setosa	50	0	0
versicolor	0	45	5
virginica	0	0	50

Colin Carroll

Network Data

Too much to say for this, here are some pictures on clustering.



Colin Carroll

Network Data

Block matrix with 20 blocks

Colin Carroll

Network Data



Colin Carroll

Network Data



Colin Carroll

Supervised learning

- Divide our data set into 120 training examples and 30 test examples.
- Not magic.
- Easiest if we assume all data is 2 dimensional.

Colin Carroll

Supervised learning

Broadly speaking, we will have a training set $\{x_j, \hat{y}_j\}$, where x_j are data and y_j are categories, and a loss function L that penalizes our model for being wrong. We use the training set to train a (hopefully simple) function f so that $\sum_j L(f(x_j), \hat{y}_j)$ is minimized.

Colin Carroll

Naive Bayes

$$p(\text{event}|\text{feature}) = \frac{p(\text{event})p(\text{feature}|\text{event})}{p(\text{feature})}$$

•
$$p(tb) = 0.005$$

(so p(positive) = 0.055).

$$p(tb|positive) = \frac{0.005 * 0.99}{0.055} = 9\%.$$

Colin Carroll

Naive Bayes

- Estimate the distribution of sepal lengths and widths for each iris species
- Q Given some sepal lengths and widths, we'll be able to calculate the probability the event belongs to each category
- Choose the category with the highest probability (*posterior*)

Called *naive* Bayes because we assume sepal length is independent from sepal width. Which isn't super realistic. Still does ok!

Colin Carroll

Naive Bayes



Colin Carroll

K nearest neighbors

Ask your K nearest neighbors which group you belong in.



Colin Carroll

K nearest neighbors

Results with 2-dimensional training data:

	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	8	2
virginica	0	4	6

Results with 4-dimensional training data:

	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	0
virginica	0	1	9

Colin Carroll

Decision tree

Find which question divides your training set the most homogeneously, repeat.



Colin Carroll

Decision tree



Colin Carroll

Decision Tree

Results with 4-dimensional training data:

	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	0
virginica	0	3	7

Colin Carroll

Support vector machines

Divide the space using a hyperplane.



Colin Carroll

Support vector machines

	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	0
virginica	0	2	8

Colin Carroll

Neural nets

Compose some activation functions together.

$$\sum_{j} \frac{1}{1 + e^{x_j \cdot w_j}}$$



Colin Carroll

Neural nets

prediction			
	setosa	versicolor	virginica
setosa	10	0	0
versicolor	0	10	0
virginica	0	3	7