TDA231 Classification: Bayes and Naive Bayes

Aristide Tossou aristide@chalmers.se

Chalmers University of Technology

April 1, 2019

Reference

The content and the slides are adapted from

S. Rogers and M. Girolami, A First Course in Machine Learning (FCML), 2nd edition, Chapman & Hall/CRC 2016, ISBN: 9781498738484

Data Representation

- Data objects e.g. email texts or images are represented by fixed dimension vectors, each dimension is called a feature.
- Traditionally (and still) the features were hand-crafted by domain experts (linguists, image researchers).
- Recently with Deep Learning, one tries to learn the features (next week!).

Classification



A set of N objects with attributes (usually vector) \mathbf{x}_n .

- Each object has an associated response (or label) t_n .
- Binary classification: $t_n = \{0, 1\}$ or $t_n = \{-1, 1\}$,
 - (depends on algorithm).
- Multi-class classification: $t_n = \{1, 2, \dots, K\}$.

Classification

- lnput is training data N pairs $(\mathbf{x}_n, t_n), n = 1 \cdots N$.
- Our algorithm should use those to produce a function f that we can apply to a new data point, a test point x_{new} to classify it.
- Binary classification: $t_n = \{0, 1\}$ or $t_n = \{-1, 1\}$,
 - (depends on algorithm).
- Multi-class classification: $t_n = \{1, 2, \dots, K\}$.

Classification syllabus

- ► 4 classification algorithms.
- Of which:
 - 2 are probabilistic.
 - Bayes classifier.
 - Logistic regression.
 - 2 are non-probabilistic.
 - K-nearest neighbours.
 - Support Vector Machines.
- There are many others!

Probabilistic v non-probabilistic classifiers

Classifier is trained on $\mathbf{x}_1, \ldots, \mathbf{x}_N$ and t_1, \ldots, t_N and then used to classify \mathbf{x}_{new} .

Probabilistic classifiers produce a probability of class membership P(t_{new} = k|x_{new}, X, t)

• e.g. binary classification: $P(t_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$ and $P(t_{\text{new}} = 0 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$.

Non-probabilistic classifiers produce a hard assignment
 e.g. t_{new} = 1 or t_{new} = 0.

Which to choose depends on application....

Probabilistic v non-probabilistic classifiers

Probabilities provide us with more information – $P(t_{\text{new}} = 1) = 0.6$ is more useful than $t_{\text{new}} = 1$.

- Tells us how sure the algorithm is.
- Particularly important where cost of misclassification is high and imbalanced.
 - e.g. Diagnosis: telling a diseased person they are healthy is much worse than telling a healthy person they are diseased.
- Extra information (probability) often comes at a cost.
- For large datasets, might have to go with non-probabilistic.

Our first probabilistic classifier is based on Bayes rule:

$$P(t_{\text{new}} = k | \mathbf{X}, \mathbf{t}, \mathbf{x}_{\text{new}})$$

=
$$\frac{P(\mathbf{x}_{\text{new}} | t_{\text{new}} = k, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = k)}{\sum_{j} p(\mathbf{x}_{\text{new}} | t_{\text{new}} = j, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = j)}$$

We need to define a likelihood and a prior and we're done!

Bayes classifier - likelihood

 $p(\mathbf{x}_{\text{new}}|t_{\text{new}} = k, \mathbf{X}, \mathbf{t})$

- How likely is x_{new} if it is in class k? (not necessarily a probability...)
- We are free to define this *class-conditional distribution* as we like.
- Will depend on type of data.
- e.g.
 - Data are *D*-dimensional vectors of real values Gaussian likelihood.
 - Data are number of heads in N coin tosses Binomial likelihood.
- In both cases, training data with t = k used to determine parameters of likelihood for class k (e.g. Gaussian mean and covariance).

Bayes classifier - prior

$$P(t_{new} = k)$$

x_{new} not present.

Used to specify prior probabilities for different classes.

▶ e.g.

There are far fewer instances of class 0 than class 1: P(t_{new} = 1) > P(t_{new} = 0).

- No prior preference: $P(t_{new} = 0) = P(t_{new} = 1)$.
- Class 0 is very rare: $P(t_{new} = 0) \ll P(t_{new} = 1)$.

Naive-Bayes

- Naive-Bayes makes the following additional likelihood assumption:
- ▶ The components of x_{new} are independent for a particular class:

$$p(\mathbf{x}_{\mathsf{new}}|t_{\mathsf{new}}=k,\mathbf{X},\mathbf{t}) = \prod_{d=1}^{D} p(x_d^{\mathsf{new}}|t_{\mathsf{new}}=k,\mathbf{X},\mathbf{t})$$

- Where D is the number of dimensions and x_d^{new} is the value of the dth one.
- Often used when D is high:
 - Fitting D uni-variate distributions is easier than fitting one D-dimensional one.

Bayes classifier, example 1



- Each object has two attributes: $\mathbf{x} = [x_1, x_2]^{\mathsf{T}}$.
- K = 3 classes.
- We'll use Gaussian class-conditional distributions (with Naive-Bayes assumption).

•
$$P(t_{new} = k) = 1/K$$
 – uniform prior.

Step 1: fitting the class-conditional densities



Step 2: Evaluate densities at test point



• Remember that we assumed $P(t_{new} = k) = 1/K$.

$$P(t_{\text{new}} = k | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}) = \frac{p(\mathbf{x}_{\text{new}} | t_{\text{new}} = k, \mathbf{X}, \mathbf{t}) p(t_{\text{new}} = k)}{\sum_{j} p(\mathbf{x}_{\text{new}} | t_{\text{new}} = j, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = j)}$$

$$P(t_{new} = 1 | \dots)$$

$$P(t_{new} = 1 | \dots)$$

• Remember that we assumed $P(t_{new} = k) = 1/K$.

$$P(t_{\text{new}} = k | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t}) = \frac{p(\mathbf{x}_{\text{new}} | t_{\text{new}} = k, \mathbf{X}, \mathbf{t}) p(t_{\text{new}} = k)}{\sum_{j} p(\mathbf{x}_{\text{new}} | t_{\text{new}} = j, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = j)}$$

$$P(t_{new} = 2 | \dots)$$

$$q_{j}$$

$$P(t_{new} = 2 | \dots)$$

$$P(t_{new} = 2 | \dots)$$

$$P(t_{new} = 2 | \dots)$$

• Remember that we assumed $P(t_{new} = k) = 1/K$.



Bayes classifier, example 2

Data are number of heads in 20 tosses (repeated 50 times for each) from one of two coins:

• Coin 1
$$(t_n = 0)$$
: $x_n = 4, 7, 7, 7, 4, ...$

• Coin 2 $(t_n = 1)$: $x_n = 18, 16, 18, 14, 17, ...$

Use binomial class conditional densities:

$$P(x_n|r_k) = \begin{pmatrix} 20 \\ x_n \end{pmatrix} r^{x_n} (1-r)^{20-x_n}$$

- Where r_k is the probability that coin k lands heads on any particular toss.
- Problem predict the coin, t_{new} given a new count, x_{new}.
- (Again assume $P(t_{new} = k) = 1/K$)

Fit the class conditionals...

Fitting is just finding *r_k*:

$$r_k = \frac{1}{20N_k} \sum_{n:t_n=k} x_n$$

▶ $r_0 = 0.287$, $r_1 = 0.706$.



$$P(t_{\text{new}} = k | x_{\text{new}}, \mathbf{X}, \mathbf{t}) = \frac{p(x_{\text{new}} | t_{\text{new}} = k, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = k)}{\sum_{j} p(x_{\text{new}} | t_{\text{new}} = j, \mathbf{X}, \mathbf{t}) P(t_{\text{new}} = j)}$$



Bayes classifier – summary

- Decision rule based on Bayes rule.
- Choose and fit class conditional densities.
- Decide on prior.
- Compute predictive probabilities.
- Naive-Bayes:
 - Assume that the dimensions of x are independent within a particular class.
 - Our Gaussian used the Naive Bayes assumption (could have written p(x|t = k,...) as product of two independent Gaussians).