

Bài tập thực hành: Phân Lớp (Classification)

Bộ môn: KHMT

Thời gian làm bài: 2 tuần (Xem deadline trong link nộp bài trên moodle)

Nộp bài:

- Nộp bài lên moodle.
- Đặt tên bài nộp theo định dạng MSSV1_MSSV2.rar. Trong đó bao gồm:
 - Tập tin báo cáo.
 - Các tập tin dữ liệu theo yêu cầu của bài tập

Các hành vi sử dụng toàn bộ/một phần bài làm của người khác sẽ bị 0 điểm cho toàn bộ phần thực hành

Bài 1. Trả lời ngắn gọn các câu hỏi sau:

- Tại sao phân lớp Bayes (Bayesian classification) được gọi là “naive” (“ngây thơ”)?
- Tại sao cần có bước tỉa nhánh (tree pruning) trong cây quyết định?
- Khi dùng cây quyết định thì có thể sẽ bị overfitting. Overfitting là gì? Và nêu một vài phương pháp để hạn chế overfitting khi dùng cây quyết định.

Bài 2. Phân Lớp.

Cho dữ liệu Mushroom (các loại nấm – xem link trên moodle). Thuộc tính cần phân lớp là nấm độc (p) hay không độc (ăn được - e).

Mô tả: có một vài giá trị không có trong file dữ liệu vì dữ liệu lớn và đã được chọn random thành tập nhỏ như bài tập.

Tên thuộc tính	Mô tả	Giá trị
Class	Lớp phân loại	e = ăn được p = có độc
Shape	Hình dáng	b = hình chuông, c = hình nón x = lồi f = phẳng k = có u s = lõm
Surface	Bề mặt	f = có xơ g = có đường rãnh (nếp) y = có vảy s = trơn
Color	Màu sắc	n = màu nâu

Khác thác dữ liệu & Ứng dụng

		b = màu da bò c = màu nâu vàng r = màu xanh lá p = màu hồng u = màu tím e = màu đỏ w = màu trắng y = màu vàng
Bruise	Có vết thâm không	t = bị thâm f = không màu
Odor	Có mùi hương	a = mùi quả hạnh l = mùi cây hồi c = mùi crê-ô-dốt y = mùi cá f = mùi thối m = mùi mốc n = không mùi p =mùi hăng s = mùi gia vị
Spore_Color	Màu bào tử (nằm)	k = màu đen n = màu nâu b = màu da bò h = màu sô-cô-la r = màu xanh lá o = màu cam u = màu tím w = màu trắng y = màu vàng
Population	Sự phân bố	s = mọc phân tán, rải rác c = mọc chùm, cụm n = mọc đông đảo, dày đặc a = mọc phong phú v = mọc nhiều y = mọc đơn độc
Habitat	Môi trường sinh sống	g = đồng cỏ l = lá cây m = đồng cỏ (và các loại cây không phải thân gỗ khác) p = đường đi u = đô thị w = bãi rác d = thân gỗ

1. Dùng weka chạy Naïve Bayes cho tập train.csv để phân lớp một loại nấm là không có độc (e) hay có độc (p).
2. Từ model của tập train đã chạy được, chạy weka cho tập test.csv. So sánh kết quả của lớp phân lớp với kết quả tập test.csv và tính xem độ chính xác của phân lớp Naïve Bayes.

Bài 3. SVM

Thực hiện các yêu cầu bài 2 với thuật toán phân lớp SVM.

Bài 4: kNN

Cài đặt thuật toán kNN với các yêu cầu sau:

- **Input:**
 - ✓ U là mẫu cần phân lớp.
 - ✓ T là tập huấn luyện: $T = (t_{1,1}, t_{1,2}, t_{1,3}, \dots, t_{1,n}), \dots, (t_{m,1}, t_{m,2}, t_{m,3}, \dots, t_{m,n})$
 - ✓ Thuộc tính $t_{i,n}$ là nhãn (label) của T_i
 - ✓ m là số lượng mẫu trong tập huấn luyện
 - ✓ n là số lượng thuộc tính trong mỗi mẫu.
 - ✓ k là số lượng láng giềng gần nhất ta cần tìm
- **Output:** Lớp của mẫu U

Bài 5: (Cộng Điểm)

1. Dùng C/C++/C# để thực hiện Naïve Bayes. Rồi so sánh kết quả đạt được và kết quả chạy từ Weka giống nhau hay khác nhau? Vì sao?
2. Thay vì Naïve Bayes, hãy dùng kNN ở bài 4 để chạy cho dữ liệu Mushroom.

Qui định:

- Làm bài theo nhóm. Mỗi nhóm tối đa 2 sinh viên.
- Hạn nộp: xem trên Moodle
- Bài nộp gồm file pdf/doc/docx trả lời câu hỏi lý thuyết, có đánh giá công việc từng cá nhân trong nhóm + thư mục source code.
- Đặt tất cả các nội dung được yêu cầu nộp trong thư mục có tên MSSV1_MSSV2, nén lại thành tập tin .zip hoặc .rar. Đại diện thay mặt nhóm để nộp ở link tương ứng trên Moodle.
- Sinh viên có thể viết chương trình bằng ngôn ngữ C/C++/C#.
- Các bài làm giống nhau hay chép code từ nơi khác sẽ bị 0 điểm.