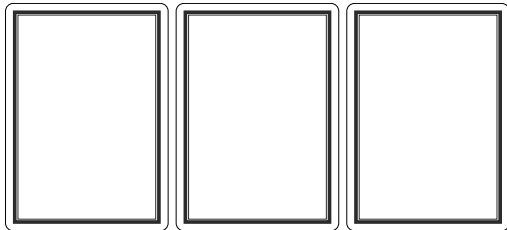


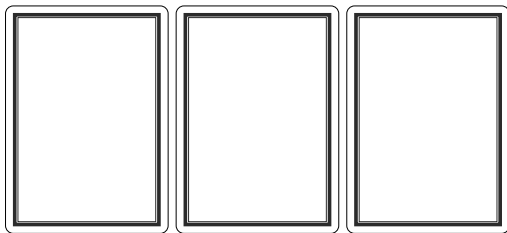


(α)

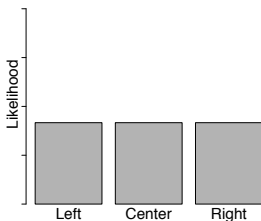


(β)

Σχήμα: Ένα παιχνίδι βρες τη ντάμα

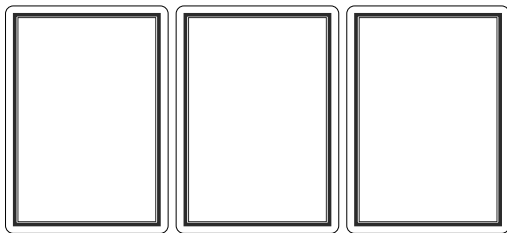


(α')

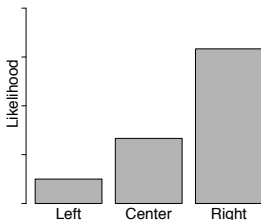


(β')

Σχήμα: Ένα παιχνίδι βρες τη ντάμα: (a) τα χαρτιά μοιρασμένα κλειστά πάνω στο τραπέζι· και (b) οι αρχικές πιθανότητες η ντάμα να καταλήξει σε κάθε θέση.

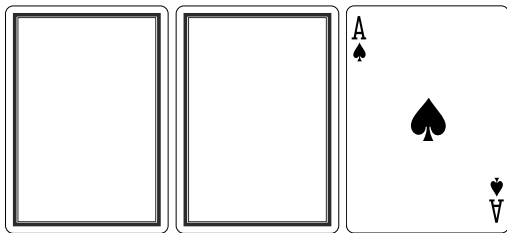


(α')

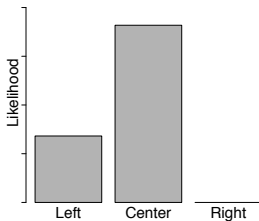


(β')

Σχήμα: Ένα παιχνίδι βρες τη ντάμα: (a) τα χαρτιά μοιρασμένα κλειστά πάνω στο τραπέζι· και (b) ένα αναθεωρημένο σύνολο πιθανοτήτων για τη θέση της ντάμας με βάση τα δεδομένα που συλλέχθηκαν.

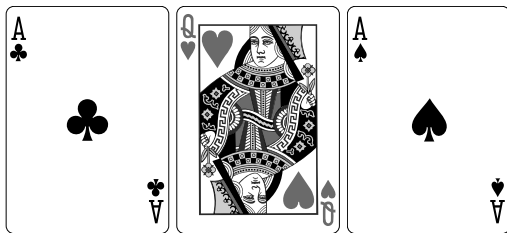


(α')



(β')

Σχήμα: Ένα παιχνίδι βρες τη ντάμα: (a) το σύνολο των χαρτιών αφού ο αέρας αναποδογυρίζει εκείνο στα δεξιά· (b) οι αναθεωρημένες πιθανότητες για τη θέση της ντάμας με βάση αυτό το νέο στοιχείο.



Σχήμα: Ένα παιχνίδι βρες τη ντάμα: οι τελικές θέσεις των χαρτιών στο παιχνίδι.

Θεμελιώδεις Έννοιες

Πίνακας: Ένα απλό σύνολο δεδομένων για τη διάγνωση MENINGITIS, με περιγραφικά χαρακτηριστικά που περιγράφουν την παρουσία ή απουσία τριών κοινών συμπτωμάτων της νόσου: HEADACHE, FEVER και VOMITING.

ID	HEADACHE	FEVER	VOMITING	MENINGITIS
1	true	true	false	false
2	false	true	false	false
3	true	false	true	false
4	true	false	true	false
5	false	true	false	true
6	true	false	true	false
7	true	false	true	false
8	true	false	true	true
9	false	true	false	false
10	true	false	true	true

Πιθανοτική ερμηνεία ενός συνόλου δεδομένων

- Κάθε χαρακτηριστικό → **τυχαία μεταβλητή**
- Όλοι οι πιθανοί συνδυασμοί τιμών → **δειγματικός χώρος**
- Κάθε γραμμή του συνόλου δεδομένων → **πείραμα**
- Κάθε συγκεκριμένη αντιστοίχιση τιμών στα περιγραφικά χαρακτηριστικά → **συμβάν**

Παράδειγμα

Κάθε γραμμή στον Πίνακα 6.1 αποτελεί ένα πείραμα και οι τιμές των περιγραφικών χαρακτηριστικών της γραμμής ορίζουν ένα διακριτό συμβάν.

Το Θεώρημα του Bayes

$$P(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

$$P(d|t) = \frac{P(t|d)P(d)}{P(t)}$$

$$\begin{aligned} P(t) &= P(t|d)P(d) + P(t|\neg d)P(\neg d) \\ &= (0.99 \times 0.0001) + (0.01 \times 0.9999) = 0.0101 \end{aligned}$$

$$\begin{aligned} P(d|t) &= \frac{0.99 \times 0.0001}{0.0101} \\ &= 0.0098 \end{aligned}$$

Για να υπολογίσουμε μια πιθανότητα χρησιμοποιώντας το Γενικευμένο Θεώρημα Bayes, πρέπει να υπολογίσουμε τρεις πιθανότητες:

- την $P(t = I)$, δηλαδή την προηγούμενη πιθανότητα (prior probability) του χαρακτηριστικού στόχου t που έχει το επίπεδο I
- την $P(q[1], \dots, q[m])$, δηλαδή τη μικτή πιθανότητα των περιγραφικών χαρακτηριστικών ενός στιγμιοτύπου ερωτήματος που παίρνει ένα συγκεκριμένο σύνολο τιμών
- την $P(q[1], \dots, q[m] \mid t = I)$, δηλαδή τη δεσμευμένη πιθανότητα των περιγραφικών χαρακτηριστικών ενός στιγμιοτύπου ερωτήματος που παίρνει ένα συγκεκριμένο επίπεδο τιμών, δεδομένου του χαρακτηριστικού-στόχου που έχει το επίπεδο I

- Μπορούμε να υπολογίσουμε τις απαιτούμενες πιθανότητες απευθείας από τα δεδομένα. Για παράδειγμα, μπορούμε να υπολογίσουμε τα $P(m)$ και $P(h, \neg f, v)$ ως εξής:

$$P(m) = \frac{|\{d_5, d_8, d_{10}\}|}{|\{d_1, d_2, d_3, d_4, d_5, d_6, d_7, d_8, d_9, d_{10}\}|} = \frac{3}{10} = 0.3$$

$$P(h, \neg f, v) = \frac{|\{d_3, d_4, d_6, d_7, d_8, d_{10}\}|}{|\{d_1, d_2, d_3, d_4, d_5, d_6, d_7, d_8, d_9, d_{10}\}|} = \frac{6}{10} = 0.6$$

- Ωστόσο, ως άσκηση θα χρησιμοποιήσουμε τον κανόνα αλυσίδας για να υπολογίσουμε:

$$P(h, \neg f, v \mid m) = ?$$

ID	HEADACHE	FEVER	VOMITING	MENINGITIS
1	true	true	false	false
2	false	true	false	false
3	true	false	true	false
4	true	false	true	false
5	false	true	false	true
6	true	false	true	false
7	true	false	true	false
8	true	false	true	true
9	false	true	false	false
10	true	false	true	true

Μπεϋζιανή Πρόβλεψη

- Χρησιμοποιώντας τον κανόνα αλυσίδας υπολογίστε:

$$\begin{aligned}P(h, \neg f, v | m) &= P(h | m) \times P(\neg f | h, m) \times P(v | \neg f, h, m) \\&= \frac{|\{\mathbf{d}_8, \mathbf{d}_{10}\}|}{|\{\mathbf{d}_5, \mathbf{d}_8, \mathbf{d}_{10}\}|} \times \frac{|\{\mathbf{d}_8, \mathbf{d}_{10}\}|}{|\{\mathbf{d}_8, \mathbf{d}_{10}\}|} \times \frac{|\{\mathbf{d}_8, \mathbf{d}_{10}\}|}{|\{\mathbf{d}_8, \mathbf{d}_{10}\}|} \\&= \frac{2}{3} \times \frac{2}{2} \times \frac{2}{2} = 0.6666\end{aligned}$$

- Ο αντίστοιχος υπολογισμός για το $P(\neg m | h, \neg f, v)$ είναι:

$$\begin{aligned} P(\neg m | h, \neg f, v) &= \frac{P(h, \neg f, v | \neg m) \times P(\neg m)}{P(h, \neg f, v)} \\ &= \frac{\left(P(h | \neg m) \times P(\neg f | h, \neg m) \right)}{P(h, \neg f, v)} \\ &\quad \times \frac{P(v | \neg f, h, \neg m) \times P(\neg m)}{P(h, \neg f, v)} \\ &= \frac{0.7143 \times 0.8 \times 1.0 \times 0.7}{0.6} = 0.6667 \end{aligned}$$

$$P(m|h, \neg f, v) = 0.3333$$

$$P(\neg m|h, \neg f, v) = 0.6667$$

- Αυτοί οι υπολογισμοί μάς λένε ότι είναι δύο φορές πιο πιθανό ο ασθενής να μην έχει μηνιγγίτιδα απ’ ό,τι να έχει, παρόλο που ο ασθενής υποφέρει από πονοκέφαλο και κάνει εμετό!

- Ο ακριβής υπολογισμός των πιθανοτήτων για κάθε ένα από τα πιθανά επίπεδα-στόχους είναι συχνά πολύ χρήσιμος για έναν άνθρωπο που λαμβάνει αποφάσεις, όπως για παράδειγμα έναν γιατρό.
- Ωστόσο, αν προσπαθούμε να κατασκευάσουμε ένα προγνωστικό μοντέλο που αναθέτει αυτόματα ένα επίπεδο-στόχο σε ένα στιγμιότυπο ερωτήματος, τότε πρέπει να αποφασίσουμε πώς το μοντέλο θα κάνει μια πρόβλεψη με βάση τις υπολογισμένες πιθανότητες.
- Ο προφανής τρόπος για να το κάνουμε αυτό είναι να ζητήσουμε από το μοντέλο να επιστρέψει το επίπεδο-στόχο που έχει την υψηλότερη εκ των υστέρων πιθανότητα, δοθείσης της κατάστασης των περιγραφικών χαρακτηριστικών στο ερώτημα.
- Ένα μοντέλο πρόβλεψης που λειτουργεί με αυτόν τον τρόπο κάνει μια μέγιστη εκ των υστέρων (maximum a posteriori (MAP)) πρόβλεψη.

ID	HEADACHE	FEVER	VOMITING	MENINGITIS
1	true	true	false	false
2	false	true	false	false
3	true	false	true	false
4	true	false	true	false
5	false	true	false	true
6	true	false	true	false
7	true	false	true	false
8	true	false	true	true
9	false	true	false	false
10	true	false	true	true

HEADACHE	FEVER	VOMITING	MENINGITIS
true	true	false	?

ID	HEADACHE	FEVER	VOMITING	MENINGITIS
1	true	true	false	false
2	false	true	false	false
3	true	false	true	false
4	true	false	true	false
5	false	true	false	true
6	true	false	true	false
7	true	false	true	false
8	true	false	true	true
9	false	true	false	false
10	true	false	true	true

$$P(m \mid h, f, \neg v) = ?$$

$$P(\neg m \mid h, f, \neg v) = ?$$

$$\begin{aligned} P(m | h, f, \neg v) &= \frac{\left(P(h|m) \times P(f | h, m) \right. \\ &\quad \left. \times P(\neg v | f, h, m) \times P(m) \right)}{P(h, f, \neg v)} \\ &= \frac{0.6666 \times 0 \times 0 \times 0.3}{0.1} = 0 \end{aligned}$$

$$\begin{aligned} P(\neg m \mid h, f, \neg v) &= \frac{\left(P(h \mid \neg m) \times P(f \mid h, \neg m) \right. \\ &\quad \left. \times P(\neg v \mid f, h, \neg m) \times P(\neg m) \right)}{P(h, f, \neg v)} \\ &= \frac{0.7143 \times 0.2 \times 1.0 \times 0.7}{0.1} = 1.0 \end{aligned}$$

$$P(m \mid h, f, \neg v) = 0$$

$$P(\neg m \mid h, f, \neg v) = 1.0$$

- Υπάρχει κάτι παράξενο σε αυτά τα αποτελέσματα!

Κατάρρα της Διαστατικότητας

Καθώς αυξάνεται ο αριθμός των περιγραφικών χαρακτηριστικών, αυξάνεται και ο αριθμός των πιθανών γεγονότων δέσμευσης. Κατά συνέπεια, απαιτείται εκθετική αύξηση του μεγέθους του συνόλου δεδομένων κάθε φορά που προστίθεται ένα νέο περιγραφικό χαρακτηριστικό, ώστε για οποιαδήποτε δεσμευμένη πιθανότητα να υπάρχουν αρκετά παραδείγματα στο σύνολο εκπαίδευσης που να ταιριάζουν στις συνθήκες, ώστε η προκύπτουσα πιθανότητα να είναι λογική.

- Η πιθανότητα ένας ασθενής που έχει πονοκέφαλο και πυρετό να έχει μηνιγγίτιδα θα πρέπει να είναι μεγαλύτερη από το μηδέν!
- Το σύνολο δεδομένων μας δεν είναι αρκετά μεγάλο → το μοντέλο μας κάνει **υπερπροσαρμογή** στα δεδομένα εκπαίδευσης.
- Οι έννοιες της **δεσμευμένης ανεξαρτησίας** και της **παραγοντοποίησης** μπορούν να μας βοηθήσουν να ξεπεράσουμε αυτό το ελάττωμα της τρέχουσας προσέγγισής μας.

- Αν η γνώση ενός γεγονότος δεν έχει καμία επίδραση στην πιθανότητα ενός άλλου γεγονότος, και *αντίστροφα*, τότε τα δύο γεγονότα είναι **ανεξάρτητα** μεταξύ τους.
- Αν δύο γεγονότα X και Y είναι ανεξάρτητα τότε:

$$P(X|Y) = P(X)$$

$$P(X, Y) = P(X) \times P(Y)$$

- Θυμηθείτε ότι όταν δύο γεγονότα είναι εξαρτημένα αυτοί οι κανόνες είναι:

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}$$

$$P(X, Y) = P(X|Y) \times P(Y) = P(Y|X) \times P(X)$$

- Η πλήρης ανεξαρτησία μεταξύ γεγονότων είναι αρκετά σπάνια.
- Ένα πιο συνηθισμένο φαινόμενο είναι δύο ή περισσότερα γεγονότα να είναι ανεξάρτητα εφόσον γνωρίζουμε ότι έχει συμβεί ένα τρίτο γεγονός.
- Αυτό είναι γνωστό ως **δεσμευμένη ανεξαρτησία**.

- Αν το γεγονός $t = l$ προκαλεί τα γεγονότα $\mathbf{q}[1], \dots, \mathbf{q}[m]$, τότε τα γεγονότα $\mathbf{q}[1], \dots, \mathbf{q}[m]$ είναι δεσμευμένα ανεξάρτητα μεταξύ τους δεδομένης της γνώσης του $t = l$, και ο ορισμός του κανόνα αλυσίδας μπορεί να απλοποιηθεί ως εξής:

$$\begin{aligned}
 P(\mathbf{q}[1], \dots, \mathbf{q}[m] \mid t = l) \\
 &= P(\mathbf{q}[1] \mid t = l) \times P(\mathbf{q}[2] \mid t = l) \times \dots \times P(\mathbf{q}[m] \mid t = l) \\
 &= \prod_{i=1}^m P(\mathbf{q}[i] \mid t = l)
 \end{aligned}$$

ID	HEADACHE	FEVER	VOMITING	MENINGITIS
1	true	true	false	false
2	false	true	false	false
3	true	false	true	false
4	true	false	true	false
5	false	true	false	true
6	true	false	true	false
7	true	false	true	false
8	true	false	true	true
9	false	true	false	false
10	true	false	true	true

- Υπολογίστε τους παράγοντες από τα δεδομένα.

$$Factor_1 : \langle P(M) \rangle$$

$$Factor_2 : \langle P(h|m), P(h|\neg m) \rangle$$

$$Factor_3 : \langle P(f|m), P(f|\neg m) \rangle$$

$$Factor_4 : \langle P(v|m), P(v|\neg m) \rangle$$

Πίνακας: Ένα σύνολο δεδομένων από πεδίο ανίχνευσης απάτης σε αιτήσεις δανείων.

ID	CREDIT HISTORY	GUARANTOR/ COAPPLICANT	ACCOMODATION	FRAUD
1	current	none	own	true
2	paid	none	own	false
3	paid	none	own	false
4	paid	guarantor	rent	true
5	arrears	none	own	false
6	arrears	none	own	true
7	current	none	own	false
8	arrears	none	own	false
9	current	none	rent	false
10	none	none	own	true
11	current	coapplicant	own	false
12	current	none	own	true
13	current	none	rent	true
14	paid	none	own	false
15	arrears	none	own	false
16	current	none	own	false
17	arrears	coapplicant	rent	false
18	arrears	none	free	false

$P(fr)$	$=$	0.3	$P(\neg fr)$	$=$	0.7
$P(CH = 'none' fr)$	$=$	0.1666	$P(CH = 'none' \neg fr)$	$=$	0
$P(CH = 'paid' fr)$	$=$	0.1666	$P(CH = 'paid' \neg fr)$	$=$	0.2857
$P(CH = 'current' fr)$	$=$	0.5	$P(CH = 'current' \neg fr)$	$=$	0.2857
$P(CH = 'arrears' fr)$	$=$	0.1666	$P(CH = 'arrears' \neg fr)$	$=$	0.4286
$P(GC = 'none' fr)$	$=$	0.8334	$P(GC = 'none' \neg fr)$	$=$	0.8571
$P(GC = 'guarantor' fr)$	$=$	0.1666	$P(GC = 'guarantor' \neg fr)$	$=$	0
$P(GC = 'coapplicant' fr)$	$=$	0	$P(GC = 'coapplicant' \neg fr)$	$=$	0.1429
$P(ACC = 'own' fr)$	$=$	0.6666	$P(ACC = 'own' \neg fr)$	$=$	0.7857
$P(ACC = 'rent' fr)$	$=$	0.3333	$P(ACC = 'rent' \neg fr)$	$=$	0.1429
$P(ACC = 'free' fr)$	$=$	0	$P(ACC = 'free' \neg fr)$	$=$	0.0714

Πίνακας: Οι πιθανότητες που χρειάζεται ένα μοντέλο πρόβλεψης Naive Bayes, υπολογισμένες από το σύνολο δεδομένων. Υπόμνημα συμβολισμού: FR=FRAUDULENT, CH=CREDIT HISTORY, GC = GUARANTOR/COAPPLICANT, ACC = ACCOMODATION, T= 'true', F= 'false'.

$P(fr)$	$=$	0.3	$P(\neg fr)$	$=$	0.7
$P(CH = 'none' fr)$	$=$	0.1666	$P(CH = 'none' \neg fr)$	$=$	0
$P(CH = 'paid' fr)$	$=$	0.1666	$P(CH = 'paid' \neg fr)$	$=$	0.2857
$P(CH = 'current' fr)$	$=$	0.5	$P(CH = 'current' \neg fr)$	$=$	0.2857
$P(CH = 'arrears' fr)$	$=$	0.1666	$P(CH = 'arrears' \neg fr)$	$=$	0.4286
$P(GC = 'none' fr)$	$=$	0.8334	$P(GC = 'none' \neg fr)$	$=$	0.8571
$P(GC = 'guarantor' fr)$	$=$	0.1666	$P(GC = 'guarantor' \neg fr)$	$=$	0
$P(GC = 'coapplicant' fr)$	$=$	0	$P(GC = 'coapplicant' \neg fr)$	$=$	0.1429
$P(ACC = 'own' fr)$	$=$	0.6666	$P(ACC = 'own' \neg fr)$	$=$	0.7857
$P(ACC = 'rent' fr)$	$=$	0.3333	$P(ACC = 'rent' \neg fr)$	$=$	0.1429
$P(ACC = 'free' fr)$	$=$	0	$P(ACC = 'free' \neg fr)$	$=$	0.0714

CREDIT HISTORY	GUARANTOR/COAPPLICANT	ACCOMODATION	FRAUDULENT
paid	none	rent	?

Το μοντέλο γενικεύει πέρα από το σύνολο δεδομένων!

ID	CREDIT HISTORY	GUARANTOR/ CoAPPLICANT	ACCOMMODATION	FRAUD
1	current	none	own	true
2	paid	none	own	false
3	paid	none	own	false
4	paid	guarantor	rent	true
5	arrears	none	own	false
6	arrears	none	own	true
7	current	none	own	false
8	arrears	none	own	false
9	current	none	rent	false
10	none	none	own	true
11	current	coapplicant	own	false
12	current	none	own	true
13	current	none	rent	true
14	paid	none	own	false
15	arrears	none	own	false
16	current	none	own	false
17	arrears	coapplicant	rent	false
18	arrears	none	free	false
19	arrears	none	own	false
20	paid	none	own	false

CREDIT HISTORY	GUARANTOR/CoAPPLICANT	ACCOMMODATION	FRAUDULENT
paid	none	rent	<i>'false'</i>

