

**Σχήμα:** Τα (a) και (b) δείχνουν δύο δέντρα αποφάσεων που είναι συνεπή με τα παραδείγματα του συνόλου δεδομένων spam. Το (c) δείχνει τη διαδρομή που ακολουθείται στο δέντρο του (a) για να γίνει πρόβλεψη για το ερώτημα: 'Υποπτες Λέξεις = 'true', Άγνωστος Αποστολέας = 'true', Περιέχει Εικόνες = 'true'.



- Εφαρμόζουμε την ίδια προσέγγιση που χρησιμοποιήσαμε στο παιχνίδι *Guess-Who* : προτιμούμε δέντρα αποφάσεων που χρησιμοποιούν λιγότερους ελέγχους (ρηχότερα δέντρα).
- Αυτό είναι ένα παράδειγμα του ξυραφιού του Occam.



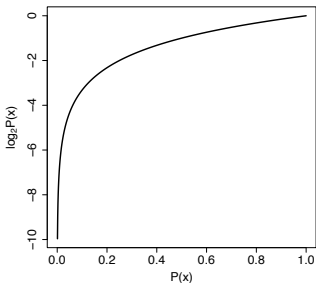
- Το μοντέλο εντροπίας του Claude Shannon ορίζει ένα υπολογιστικό μέτρο της ακαθαρσίας των στοιχείων ενός συνόλου.
- Ένας εύκολος τρόπος να κατανοήσουμε την εντροπία ενός συνόλου είναι να σκεφτούμε την αβεβαιότητα που σχετίζεται με την πρόβλεψη του αποτελέσματος όταν κάνουμε μια τυχαία επιλογή από το σύνολο.

- Η εντροπία σχετίζεται με την πιθανότητα ενός αποτελέσματος.
  - Υψηλή πιθανότητα → Χαμηλή εντροπία
  - Χαμηλή πιθανότητα → Υψηλή εντροπία
- Αν πάρουμε τον **λογάριθμο** μιας πιθανότητας και τον πολλαπλασιάσουμε με  $-1$ , παίρνουμε αυτή την αντιστοίχιση!

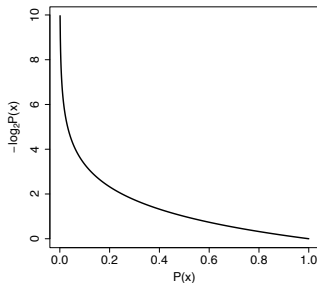
## Τι είναι ο λογάριθμος;

Θυμηθείτε ότι ο λογάριθμος του  $a$  με βάση το  $b$  είναι ο αριθμός στον οποίο πρέπει να υψώσουμε το  $b$  ώστε να πάρουμε το  $a$ .

- $\log_2(0.5) = -1$  επειδή  $2^{-1} = 0.5$
- $\log_2(1) = 0$  επειδή  $2^0 = 1$
- $\log_2(8) = 3$  επειδή  $2^3 = 8$
- $\log_5(25) = 2$  επειδή  $5^2 = 25$
- $\log_5(32) = 2.153$  επειδή  $5^{2.153} = 32$



(α')



(β')

**Σχήμα:** (α) Γράφημα που δείχνει πώς αλλάζει η τιμή του δυαδικού λογαρίθμου (λογάριθμος με βάση το 2) μιας πιθανότητας σε όλο το εύρος των τιμών πιθανότητας. (β) Η επίδραση του πολλαπλασιασμού αυτών των τιμών με  $-1$ .

- Το μοντέλο εντροπίας του Shannon είναι ένα σταθμισμένο άθροισμα των λογαρίθμων των πιθανοτήτων κάθε πιθανού αποτελέσματος όταν κάνουμε μια τυχαία επιλογή από ένα σύνολο.

$$H(t) = - \sum_{i=1}^l (P(t = i) \times \log_s(P(t = i))) \tag{1}$$

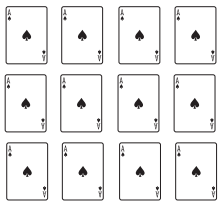
- Ποια είναι η εντροπία ενός συνόλου από 52 διαφορετικά τραπουλόχαρτα;

$$\begin{aligned}
 H(card) &= - \sum_{i=1}^{52} P(card = i) \times \log_2(P(card = i)) \\
 &= - \sum_{i=1}^{52} 0.019 \times \log_2(0.019) = - \sum_{i=1}^{52} -0.1096 \\
 &= 5.700 \text{ bits}
 \end{aligned}$$

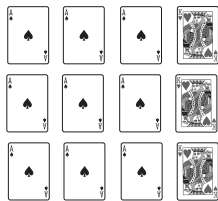
- Ποια είναι η εντροπία ενός συνόλου από 52 τραπουλόχαρτα αν τα διακρίνουμε μόνο με βάση το χρώμα τους {♥, ♣, ♦, ♠};

$$\begin{aligned}
 H(\text{suit}) &= - \sum_{l \in \{\heartsuit, \clubsuit, \diamondsuit, \spadesuit\}} P(\text{suit} = l) \times \log_2(P(\text{suit} = l)) \\
 &= - \left( (P(\heartsuit) \times \log_2(P(\heartsuit))) + (P(\clubsuit) \times \log_2(P(\clubsuit))) \right. \\
 &\quad \left. + (P(\diamondsuit) \times \log_2(P(\diamondsuit))) + (P(\spadesuit) \times \log_2(P(\spadesuit))) \right) \\
 &= - \left( \left( \frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) + \left( \frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) \right. \\
 &\quad \left. + \left( \frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) + \left( \frac{13}{52} \times \log_2\left(\frac{13}{52}\right) \right) \right) \\
 &= - \left( (0.25 \times -2) + (0.25 \times -2) \right. \\
 &\quad \left. + (0.25 \times -2) + (0.25 \times -2) \right) \\
 &= 2 \text{ bits}
 \end{aligned}$$

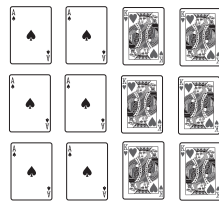
Μοντέλο Εντροπίας του Shannon



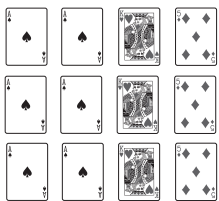
(α')  $H(card) = 0.00$



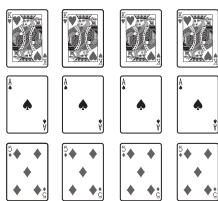
(β')  $H(card) = 0.81$



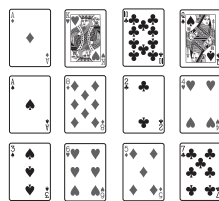
(γ')  $H(card) = 1.00$



(δ')  $H(card) = 1.50$



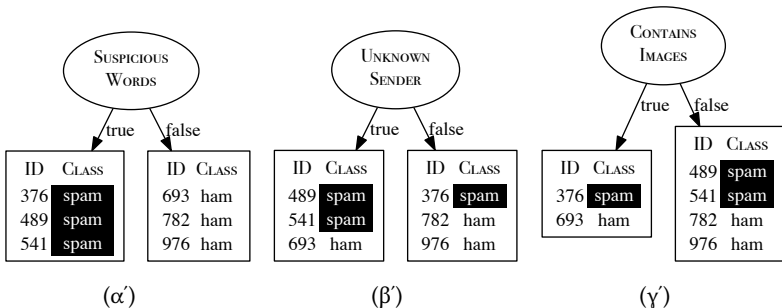
(ε')  $H(card) = 1.58$



(ς')  $H(card) = 3.58$

**Σχήμα:** Η εντροπία διαφορετικών συνόλων τραπουλόχαρτων, μετρημένη σε bits.





**Σχήμα:** Πώς χωρίζονται τα παραδείγματα του συνόλου δεδομένων spam όταν τα διαμερίζουμε με βάση καθένα από τα περιγραφικά χαρακτηριστικά του συνόλου δεδομένων spam στον Πίνακα 1 <sup>[15]</sup>

- Η διαίσθησή μας είναι ότι το ιδανικό διακριτικό χαρακτηριστικό θα διαμερίζει τα δεδομένα σε **καθαρά** υποσύνολα όπου όλα τα παραδείγματα σε κάθε υποσύνολο έχουν την ίδια κατηγορία.
  - Ύποπτες Λέξεις: τέλειος διαχωρισμός.
  - Άγνωστος Αποστολέας: μίγμα, αλλά με κάποια πληροφορία (όταν είναι 'true', τα περισσότερα παραδείγματα είναι 'spam').
  - Περιέχει Εικόνες: καμία πληροφορία.
- Ένας τρόπος να υλοποιήσουμε αυτή την ιδέα είναι να χρησιμοποιήσουμε ένα μέτρο που λέγεται **κέρδος πληροφορίας**.

## Κέρδος Πληροφορίας

- Το κέρδος πληροφορίας ενός περιγραφικού χαρακτηριστικού μπορεί να ερμηνευτεί ως μέτρο της μείωσης της συνολικής εντροπίας ενός έργου πρόβλεψης όταν ελέγχουμε αυτό το χαρακτηριστικό.

Ο υπολογισμός του κέρδους πληροφορίας περιλαμβάνει τις ακόλουθες 3 εξισώσεις:

$$H(t, \mathcal{D}) = - \sum_{l \in \text{levels}(t)} (P(t = l) \times \log_2(P(t = l))) \quad (2)$$

$$\text{rem}(d, \mathcal{D}) = \sum_{l \in \text{levels}(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{στάθμιση}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\substack{\text{εντροπία του} \\ \text{διαμερίσματος } \mathcal{D}_{d=l}}} \quad (3)$$

$$IG(d, \mathcal{D}) = H(t, \mathcal{D}) - \text{rem}(d, \mathcal{D}) \quad (4)$$

- Ως παράδειγμα, θα υπολογίσουμε το κέρδος πληροφορίας για καθένα από τα περιγραφικά χαρακτηριστικά στο σύνολο δεδομένων spam email.

- Υπολογίστε την **εντροπία** για το χαρακτηριστικό-στόχο στο σύνολο δεδομένων.

$$H(t, \mathcal{D}) = - \sum_{l \in \text{levels}(t)} (P(t = l) \times \log_2(P(t = l)))$$

ID	Ύποπτες Λέξεις	Άγνωστος Αποστολέας	Περιέχει Εικόνες	Κλάση
376	true	false	true	spam
489	true	true	false	spam
541	true	true	false	spam
693	false	true	true	ham
782	false	false	false	ham
976	false	false	false	ham

$$\begin{aligned}
 H(t, \mathcal{D}) &= - \sum_{l \in \{ 'spam', 'ham' \}} (P(t = l) \times \log_2(P(t = l))) \\
 &= - ((P(t = 'spam') \times \log_2(P(t = 'spam'))) \\
 &\quad + (P(t = 'ham') \times \log_2(P(t = 'ham')))) \\
 &= -((\frac{3}{6} \times \log_2(\frac{3}{6})) + (\frac{3}{6} \times \log_2(\frac{3}{6}))) \\
 &= 1 \text{ bit}
 \end{aligned}$$

- Υπολογίστε το **υπόλοιπο** για το χαρακτηριστικό Έγχοπτες Λέξεις στο σύνολο δεδομένων.

$$rem(d, \mathcal{D}) = \sum_{l \in levels(d)} \underbrace{\frac{|\mathcal{D}_{d=l}|}{|\mathcal{D}|}}_{\text{στάθμιση}} \times \underbrace{H(t, \mathcal{D}_{d=l})}_{\text{εντροπία του διαμερίσματος } \mathcal{D}_{d=l}}$$

$$\begin{aligned}
 &rem(\text{WORDS}, \mathcal{D}) \\
 &= \left( \frac{|\mathcal{D}_{\text{WORDS}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{WORDS}=T}) \right) + \left( \frac{|\mathcal{D}_{\text{WORDS}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{WORDS}=F}) \right) \\
 &= \left( \frac{3}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &+ \left( \frac{3}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &= \left( \frac{3}{6} \times \left( - \left( \left( \frac{3}{3} \times \log_2\left(\frac{3}{3}\right) \right) + \left( \frac{0}{3} \times \log_2\left(\frac{0}{3}\right) \right) \right) \right) \right) \\
 &+ \left( \frac{3}{6} \times \left( - \left( \left( \frac{0}{3} \times \log_2\left(\frac{0}{3}\right) \right) + \left( \frac{3}{3} \times \log_2\left(\frac{3}{3}\right) \right) \right) \right) \right) = 0 \text{ bits}
 \end{aligned}$$

- Υπολογίστε το **υπόλοιπο** για το χαρακτηριστικό Άγνωστος Αποστολέας στο σύνολο δεδομένων.

$$\begin{aligned}
 &rem(\text{SENDER}, \mathcal{D}) \\
 &= \left( \frac{|\mathcal{D}_{\text{SENDER}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{SENDER}=T}) \right) + \left( \frac{|\mathcal{D}_{\text{SENDER}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{SENDER}=F}) \right) \\
 &= \left( \frac{3}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &+ \left( \frac{3}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &= \left( \frac{3}{6} \times \left( - \left( \left( \frac{2}{3} \times \log_2\left(\frac{2}{3}\right) \right) + \left( \frac{1}{3} \times \log_2\left(\frac{1}{3}\right) \right) \right) \right) \right) \\
 &+ \left( \frac{3}{6} \times \left( - \left( \left( \frac{1}{3} \times \log_2\left(\frac{1}{3}\right) \right) + \left( \frac{2}{3} \times \log_2\left(\frac{2}{3}\right) \right) \right) \right) \right) = 0.9183 \text{ bits}
 \end{aligned}$$

- Υπολογίστε το **υπόλοιπο** για το χαρακτηριστικό Περιέχει Εικόνες στο σύνολο δεδομένων.

$$\begin{aligned}
 &rem(\text{IMAGES}, \mathcal{D}) \\
 &= \left( \frac{|\mathcal{D}_{\text{IMAGES}=T}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{IMAGES}=T}) \right) + \left( \frac{|\mathcal{D}_{\text{IMAGES}=F}|}{|\mathcal{D}|} \times H(t, \mathcal{D}_{\text{IMAGES}=F}) \right) \\
 &= \left( \frac{2}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &+ \left( \frac{4}{6} \times \left( - \sum_{l \in \{ 'spam', 'ham' \}} P(t=l) \times \log_2(P(t=l)) \right) \right) \\
 &= \left( \frac{2}{6} \times \left( - \left( \left( \frac{1}{2} \times \log_2\left(\frac{1}{2}\right) \right) + \left( \frac{1}{2} \times \log_2\left(\frac{1}{2}\right) \right) \right) \right) \right) \\
 &+ \left( \frac{4}{6} \times \left( - \left( \left( \frac{2}{4} \times \log_2\left(\frac{2}{4}\right) \right) + \left( \frac{2}{4} \times \log_2\left(\frac{2}{4}\right) \right) \right) \right) \right) = 1 \text{ bit}
 \end{aligned}$$



