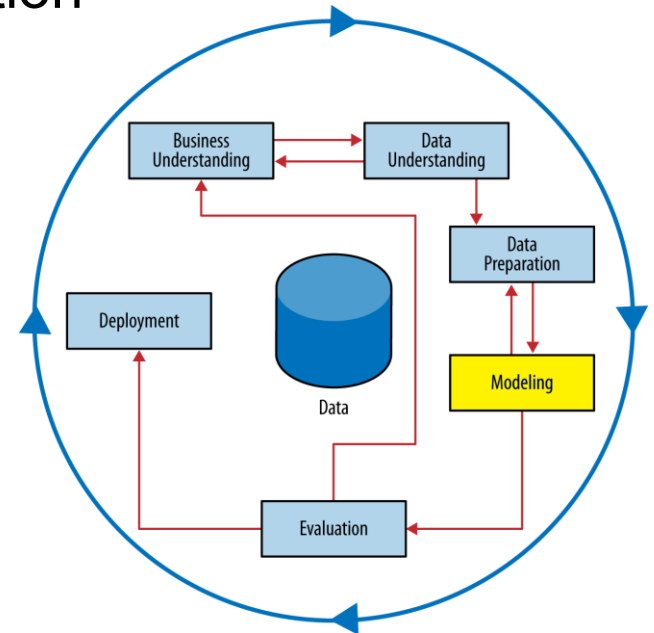


Evidence and Probabilities

PF9

Learning Goals

- ▶ Explicit evidence combination with **Bayes' Rule**
 - ▶ Take **features** as evidences, while the **classification result** as conditional result for the evidence(s).
- ▶ Probabilistic reasoning via the assumption of **conditional independence** among multiple evidences.
 - ▶ Naive Bayes
- ▶ Techniques:
 - ▶ Naive Bayes Classification
 - ▶ Evidence Lift



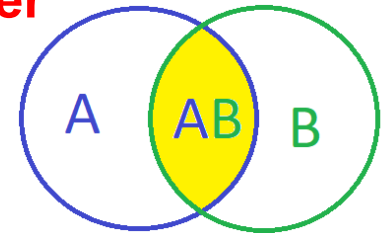
Example: An Upscale Hotel Chain

- ▶ **Target**: whether a consumer will book a room in our hotel?
- ▶ **Features**: whether the customer (1) visited our website; (2) did online shopping; (3) bought a travel book on Amazon or not.
- ▶ **Probabilistic reasoning**: combine multiple evidences probabilistically.
 - ▶ A feature as an evidence, all evidences as a feature vector $E = [e_1, \dots, e_k]$.
 - ▶ The customer has visited our website.
 - ▶ The customer has been doing online shopping.
 - ▶ The customer has bought a travel book from Amazon.
 - ▶ What is the probability that a customer will book a room (c_i) given the feature values (E) about him/her? $P(c_i|E)$
- ▶ **Evidence Lift**: how would an evidence (e_i) increase the estimated probability of booking $P(c_i|E)$, in comparison to the baseline?
 - ▶ The baseline (prior) probability considers no evidence: $p(c_i)$

Conditional and Joint Probability

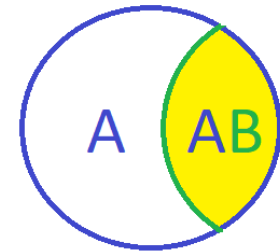
- ▶ The area AB: **joint probability A & B occur together**

$$p(A \cap B) \text{ or } p(AB)$$



- ▶ Ratio of area AB to A: **conditional probability of B given A**

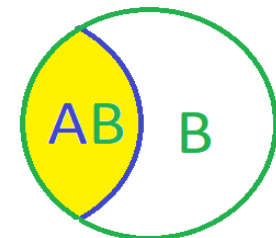
$$p(B|A) = \frac{p(AB)}{p(A)}$$



- ▶ The joint probability $p(AB) = p(B|A) \times p(A)$

- ▶ Ratio of area AB to B: **conditional probability of A given B**

$$p(A|B) = \frac{p(AB)}{p(B)}$$

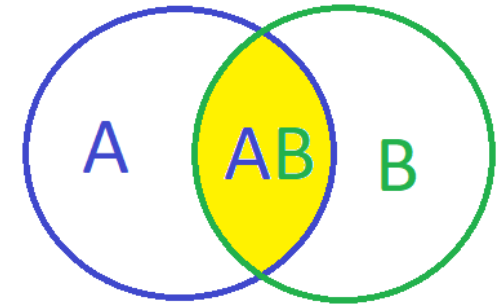


- ▶ The joint probability $p(AB) = p(A|B) \times p(B)$

$$p(AB) = p(B|A) \times p(A) = p(A|B) \times p(B)$$

Example: Conditional and Joint Probability

- ▶ Assume there are 5 days per week:
 - ▶ Student **A** comes to school 4 days every week;
 - ▶ Student **B** comes to school 3 days every week;
 - ▶ Student **A** and **B** come to school together 2 days every week.



- ▶ Please answer below questions:
 - ▶ What is the probability that **A** come today? $\frac{4}{5}$
 - ▶ What is the probability that **B** come today? $\frac{3}{5}$
 - ▶ What is the probability that **A** and **B** come together? $\frac{2}{5}$
 - ▶ What is the probability **B** would come if **A** is here? $\frac{2}{4}$
 - ▶ What is the probability **A** would come if **B** is here? $\frac{2}{3}$
- Unconditional probability (Prior probability)**
- Joint probability**
- Conditional probability**

Exercise: Gender vs. Hair Color

- ▶ Unconditional probability:

- ▶ $P(\text{Red}) = \frac{16}{30}$

- ▶ $P(\text{Yellow}) = \frac{14}{30}$

- ▶ $P(\text{Male}) = \frac{18}{30}$

- ▶ $P(\text{Female}) = \frac{12}{30}$

| | Yellow Hair | Red Hair | Total |
|--------|-------------|----------|-------|
| Male | 8 | 10 | 18 |
| Female | 6 | 6 | 12 |
| Total | 14 | 16 | 30 |

- ▶ Compute the probabilities below:

- ▶ $P(\text{Male} \cap \text{Yellow})$: joint probability to see a male person with yellow hair
 - ▶ $P(\text{Yellow} \mid \text{Male})$: conditional probability of yellow hair given he is a male
 - ▶ $P(\text{Male} \mid \text{Yellow})$: conditional probability of being male given s/he has yellow hair

Independence

- ▶ Independence is a special case of conditional probability.
 - ▶ If **A** and **B** are independent, knowing **A** tells nothing about **B**.

- ▶ If student **A** and **B** are independent:

- ▶ Conditional probability to see A given B is

$$p(A | B) = p(A) = 80\%$$

Student A comes to school
4 days every week

- ▶ Conditional probability to see B given A is

$$p(B | A) = p(B) = 60\%$$

Student B comes to school
3 days every week

- ▶ Joint probability to see A and B together should be

$$p(AB) = p(A) \times p(B) = 80\% \times 60\% = 48\%$$

Higher than the joint
probability (40%) observed

- ▶ Joint probability can also be rewritten as:

- ▶ $p(AB) = p(A) \times p(B) = p(A|B) \times p(B)$

- ▶ $p(AB) = p(B) \times p(A) = p(B|A) \times p(A)$

$$p(AB) = p(B|A) \times p(A) = p(A|B) \times p(B)$$

Bayes' Rule for Classification

- ▶ **Bayes' Rule:**

$$p(AB) = p(B|A) \times p(A) = p(A|B) \times p(B)$$

- ▶ **Conditional Probability of B on A:**

$$p(B|A) = \frac{p(A|B) \times p(B)}{p(A)}$$

- ▶ Replace A with **E** (i.e., evidence), B with **c_i** (i.e., a classification result):
 - ▶ Usually evidences **E** are easily observed, while the classification result **c_i** is something we are interested but NOT easily observed.

$$p(c_i|E) = \frac{p(E|c_i) \times p(c_i)}{p(E)}$$

Estimate Class Probability With Bayes' Rule

- ▶ Unconditional probability:

- ▶ $P(\text{Red}) = \frac{16}{30}$

- ▶ $P(\text{Yellow}) = \frac{14}{30}$

- ▶ $P(\text{Male}) = \frac{18}{30}$

- ▶ $P(\text{Female}) = \frac{12}{30}$

| | Yellow Hair | Red Hair | Total |
|--------|-------------|----------|-------|
| Male | 8 | 10 | 18 |
| Female | 6 | 6 | 12 |
| Total | 14 | 16 | 30 |

- ▶ The probability someone is male (c_i) given s/he has yellow hair (E) :

$$\begin{aligned} P(\text{Male}|\text{Yellow}) &= \frac{P(\text{Yellow}|\text{Male}) \times P(\text{Male})}{P(\text{Yellow})} \\ &= \frac{\frac{8}{18} \times \frac{18}{30}}{\frac{14}{30}} = \frac{8}{14} \end{aligned}$$

Compute $p(E)$ can be difficult.

Example 1 – Compute $p(E)$

- ▶ There are 3 coins with the same look and weight.
 - ▶ 2 coins are fair: $P(\text{Head} \mid F) = 0.5$
 - ▶ 1 coin is magic: $P(\text{Head} \mid M) = 0.75$
- ▶ One coin is picked randomly. How likely this is a **magic coin**?
 - ▶ $P(M)$: the unconditional/prior probability – the baseline

$$P(M) = 1/3$$

- ▶ Now, we throw it once: we see a **head**. What is the probability that this is a **magic coin**?
 - ▶ Evidence: Head
 - ▶ Outcome: Magic

$$P(M \mid \text{Head})$$

Example 1 – Compute $p(E)$

- ▶ Unconditional probability:
 - ▶ $P(M) = 1/3$; $P(F) = 2/3$
- ▶ Conditional probability: likelihood to see a **Head** given a magic/fair coin:
 - ▶ $P(\text{Head}|M) = 3/4$; $P(\text{Head}|F) = 1/2$
- ▶ According to Bayes' rule, the probability of being a magic coin given a head is observed:

$$P(M|\text{Head}) = \frac{P(\text{Head}|M) \times P(M)}{P(\text{Head})}$$

- ▶
$$P(\text{Head}) = P(M) \times P(\text{Head}|M) + P(F) \times P(\text{Head}|F)$$
$$= 1/3 \times 3/4 + 2/3 \times 1/2 = 7/12$$

- ▶
$$P(M|\text{Head}) = \frac{P(\text{Head}|M) \times P(M)}{P(\text{Head})}$$
$$= (3/4 \times 1/3) / (7/12) = 3/7$$

More likely this is a magic coin
if a head is observed
($3/7 > 1/3$)

Example 2 – Compute $p(E)$

- ▶ There are 3 coins with the same look and weight.
 - ▶ 2 coins are fair: $P(\text{Head} \mid F) = 0.5$
 - ▶ 1 coin is magic: $P(\text{Head} \mid M) = 0.75$
- ▶ One coin is picked randomly, and we throw it *twice*. The outcome is {Head, Head}. What is the probability that this is a magic coin?
 - ▶ Evidence: {Head, Head}
 - ▶ Outcome: Magic

$$P(M \mid HH)$$

Recap: Bayes' Rule for Classification

$$p(c_i|E) = \frac{p(E|c_i) \times p(c_i)}{p(E)}$$

- ▶ c_i is a class label (i.e., one target value).
- ▶ $E = [e_1, e_2, \dots, e_k]$ is a k -dimensional feature vector, with e_i is a feature value.
- ▶ $p(c_i|E)$: the (posterior) **conditional probability** of class c_i .
 - ▶ The probability an instance belongs to class c_i given evidence E observed.
- ▶ $p(c_i)$: the (prior) **unconditional probability** of class c_i (i.e., baseline).
 - ▶ The probability an instance is classified as c_i without any evidence.
 - ▶ Usually inferred from data: i.e., proportion of class c_i in training data.
- ▶ $p(E|c_i)$: the **conditional probability** to see evidence E in class c_i .
 - ▶ i.e., proportion of instances in class c_i with evidence E .

How can we handle multiple features/evidences?

Naive Bayes: Conditional Independence

- ▶ According to the formula of joint probability:

$$p(AB) = p(B|A) \times p(A)$$

- ▶ If there are three events:

$$p(ABC) = p(B|AC) \times p(A|C) \times p(C)$$

$$p(AB|C) = \frac{p(ABC)}{p(C)} = p(B|AC) \times p(A|C)$$

- ▶ If we assume A and B are **conditionally independent** given C, then

$$p(AB | C) = p(A | C) \times p(B | C)$$

- ▶ **Naive Bayes** assumes conditional independence among features given class c_i .

- ▶ $p(E|c_i) = p(e_1, e_2, \dots, e_k|c_i) = p(e_1|c_i) \times p(e_2|c_i) \times \dots \times p(e_k|c_i)$

- ▶ $p(e_i|c_i)$ can be estimated directly from the data.

Naive Bayes for Classification

▶ The Naive Bayes Equation:

$$p(c_i | E) = \frac{p(E | c_i) \times p(c_i)}{p(E)}$$

$$\text{▶ } p(E | c_i) = p(e_1 | c_i) \times p(e_2 | c_i) \times \dots \times p(e_k | c_i)$$

$$\begin{aligned} \text{▶ } p(E) &= p(E \cap c_0) + p(E \cap c_1) && \text{Assume binary classification} \\ &= p(E | c_0) \times p(c_0) + p(E | c_1) \times p(c_1) \\ &= p(e_1 | c_0) \times \dots \times p(e_k | c_0) \times p(c_0) + p(e_1 | c_1) \times \dots \times p(e_k | c_1) \times p(c_1) \end{aligned}$$

- ▶ Note 1: if we aim for classification, no need to compute $p(E)$ as it is the same for all classes.
- ▶ Note 2: with multiple features, $p(E)$ based on naive bayes rule is NOT equal to the proportion of instances with feature E in the data. Same logic applies to $p(E | c_i)$ and $p(c_i | E)$.

Naive Bayes Classifier

- ▶ 3 features:
 - ▶ Height = s, m, t
 - ▶ Weight = h, l, n
 - ▶ Long_hair = y, n

- ▶ Target variable:

- ▶ Gender = F, M

- ▶ For a person who is tall, normal weight, long hair, what is the predicted gender?


| ID | Height | Weight | Long_hair | Gender |
|----|--------|--------|-----------|--------|
| 1 | m | n | n | M |
| 2 | s | l | y | F |
| 3 | t | h | n | M |
| 4 | s | n | y | F |
| 5 | t | n | y | F |
| 6 | s | l | n | F |
| 7 | s | h | y | M |
| 8 | m | n | n | F |
| 9 | m | l | y | F |
| 10 | t | n | n | M |

$$p(G = F | H = t, W = n, L = y)$$

Is it 100%?

Naive Bayes Classifier

- ▶ Compute the likelihood of being **Male** given the person is **tall**, **normal weight** and **long hair**.
 - ▶ The probability of $G = M$ given the three feature values is:

Bayes' Rule 

$$P(G = M | H = t, W = n, L = y)$$
$$= \frac{P(H = t, W = n, L = y | G = M) \times P(G = M)}{P(H = t, W = n, L = y)}$$

Conditional Independence assumption 

$$= \frac{[P(H = t | G = M) \times P(W = n | G = M) \times P(L = y | G = M)] \times P(G = M)}{P(H = t, W = n, L = y)}$$

Naive Bayes Classifier

- ▶ Compute the likelihood of being **Female** given the person is **tall**, **normal weight** and **long hair**.
 - ▶ The probability of $G = F$ given the three feature values is:

$$P(G = F | H = t, W = n, L = y)$$

Bayes' Rule



$$= \frac{P(H = t, W = n, L = y | G = F) \times P(G = F)}{P(H = t, W = n, L = y)}$$

Conditional Independence assumption



$$= \frac{[P(H = t | G = F) \times P(W = n | G = F) \times P(L = y | G = F)] \times P(G = F)}{P(H = t, W = n, L = y)}$$

Naive Bayes Classifier

- Conditional probabilities of being tall:

- $$P(H = t|G = M) = \frac{2}{4}$$
- $$P(H = t|G = F) = \frac{1}{6}$$

- Conditional probabilities of having normal weight:

- $$P(W = n|G = M) = \frac{2}{4}$$
- $$P(W = n|G = F) = \frac{3}{6}$$

- Conditional probabilities of having long hair:

- $$P(L = y|G = M) = \frac{1}{4}$$
- $$P(L = y|G = F) = \frac{4}{6}$$

| ID | Height | Weight | Long_hair | Gender |
|-----------|----------|----------|-----------|----------|
| 1 | <i>m</i> | <i>n</i> | <i>n</i> | <i>M</i> |
| 2 | <i>s</i> | <i>l</i> | <i>y</i> | <i>F</i> |
| 3 | <i>t</i> | <i>h</i> | <i>n</i> | <i>M</i> |
| 4 | <i>s</i> | <i>n</i> | <i>y</i> | <i>F</i> |
| 5 | <i>t</i> | <i>n</i> | <i>y</i> | <i>F</i> |
| 6 | <i>s</i> | <i>l</i> | <i>n</i> | <i>F</i> |
| 7 | <i>s</i> | <i>h</i> | <i>y</i> | <i>M</i> |
| 8 | <i>m</i> | <i>n</i> | <i>n</i> | <i>F</i> |
| 9 | <i>m</i> | <i>l</i> | <i>y</i> | <i>F</i> |
| 10 | <i>t</i> | <i>n</i> | <i>n</i> | <i>M</i> |

Naive Bayes Classifier

- ▶ The probability of $G = M$ given the three features is:

$$\begin{aligned} & P(G = M | H = t, W = n, L = y) \\ &= \frac{[P(H = t | G = M) \times P(W = n | G = M) \times P(L = y | G = M)] \times P(G = M)}{P(H = t, W = n, L = y)} \\ &= \frac{\frac{2}{4} \times \frac{2}{4} \times \frac{1}{4} \times \frac{4}{10}}{P(H = t, W = n, L = y)} \\ &= \frac{\frac{1}{40}}{P(H = t, W = n, L = y)} \end{aligned}$$

- ▶ As $p(E)$ is the same for both class, no need to compute it if classification is the goal.

Naive Bayes Classifier

- ▶ The probability of $G = F$ given the three features is:

$$\begin{aligned} & P(G = F | H = t, W = n, L = y) \\ &= \frac{[P(H = t | G = F) \times P(W = n | G = F) \times P(L = y | G = F)] \times P(G = F)}{P(H = t, W = n, L = y)} \\ &= \frac{\frac{1}{6} \times \frac{3}{6} \times \frac{4}{6} \times \frac{6}{10}}{P(H = t, W = n, L = y)} \\ &= \frac{\frac{1}{30}}{P(H = t, W = n, L = y)} \end{aligned}$$

$$1/30 > 1/40$$

- ▶ **Conclusion:** the gender of this person (tall, normal weight, long hair) is predicted as as **Female**.
 - ▶ $p(E)$ is the same for both class, no need to compute it.


Estimate Class Probability

- ▶ Let's look at the denominator $p(E)$:

$$P(H = t, W = n, L = y)$$

$$= P(H = t, W = n, L = y | G = \textcolor{red}{M}) \times P(G = \textcolor{red}{M}) + P(H = t, W = n, L = y | G = \textcolor{green}{F}) \times P(G = \textcolor{green}{F})$$

Conditional Independence assumption


$$= [P(H = t | G = \textcolor{red}{M}) \times P(W = n | G = \textcolor{red}{M}) \times P(L = y | G = \textcolor{red}{M})] \times P(G = \textcolor{red}{M}) + [P(H = t | G = \textcolor{green}{F}) \times P(W = n | G = \textcolor{green}{F}) \times P(L = y | G = \textcolor{green}{F})] \times P(G = \textcolor{green}{F})$$

$$= \left(\frac{2}{4} \times \frac{2}{4} \times \frac{1}{4}\right) \times \frac{4}{10} + \left(\frac{1}{6} \times \frac{3}{6} \times \frac{4}{6}\right) \times \frac{6}{10} = \frac{\textcolor{red}{7}}{\textcolor{red}{120}}$$

- ▶ Then

$$\textcolor{teal}{P}(G = \textcolor{red}{M} | H = t, W = n, L = y) = \frac{\frac{1}{40}}{\textcolor{red}{P}(H=t, W=n, L=y)} = \frac{\frac{1}{40}}{\frac{7}{120}} = \frac{3}{7}$$

$$\textcolor{teal}{P}(G = \textcolor{green}{F} | H = t, W = n, L = y) = \frac{\frac{1}{30}}{\textcolor{red}{P}(H=t, W=n, L=y)} = \frac{\frac{1}{30}}{\frac{7}{120}} = \frac{4}{7}$$

Naive Bayes Classifier

► Advantages

- A very simple classifier: store the counts of classes and feature as each instance is seen.
 - $p(c_i)$: proportion of instances of class c_i in the dataset.
 - $p(e_i|c_i)$: proportion of instances in class c_i with feature value e_i .
- **Incremental learning** allows the model to learn from new data without forgetting previously learnt information.
 - Useful when data is continuously generated (e.g., streaming platforms) or too large to be processed all at once.

► Disadvantages

- Class probability estimation may not be accurate, ranking is fine.
 - Probability is overestimated for correct class but underestimated for incorrect class, as the **assumption of conditional independence** is usually unrealistic.

Gaussian Naive Bayes

- ▶ Gaussian NB works for data with **continuous features**.
 - ▶ e.g., petal width, petal length for the iris data
- ▶ Estimate the prior probability for each class c_i . $p(c_i)$
 - ▶ i.e., proportion of instances of class c_i in the training set.
- ▶ Compute each feature's **mean** and **standard deviation** in each class c_i , assuming **normal distribution**.
- ▶ When making prediction for a new instance X : $p(e_i|c_i)$
 - ▶ compute the conditional probability for X takes on feature value e_i in class c_i , using **probability density function** of continuous random variables.
- ▶ Estimate class probability for X
$$\frac{(e_1|c_i) \times \cdots \times p(e_k|c_i) \times p(c_i)}{p(E)}$$

Evidence Lift

- Assuming **unconditional independence** for all features:

- $p(E) = p(e_1) \times p(e_2) \times \cdots \times p(e_k)$

- $$p(c_i|E) = \frac{[p(e_1|c_i) \times p(e_2|c_i) \times \cdots \times p(e_k|c_i)] \times p(c_i)}{p(E)} = \frac{[p(e_1|c_i) \times p(e_2|c_i) \times \cdots \times p(e_k|c_i)] \times p(c_i)}{p(e_1) \times p(e_2) \times \cdots \times p(e_k)}$$

Naive-Naive Bayes

- Evidence lift** for a feature value (e_i) is:

$$lift_{c_i}(e_i) = \frac{p(e_i|c_i)}{p(e_i)}, \quad i = 1, \dots, k$$

$$p(c_i|E) = p(c_i) \times lift_{c_i}(e_1) \times \cdots \times lift_{c_i}(e_k)$$

- Each feature value raises/lowers the prior probability $p(c_i)$ by its lift score.
 - If $lift_{c_i}(e_i) > 1$, then $p(c_i|E) > p(c_i)$.
 - If $lift_{c_i}(e_i) < 1$, then $p(c_i|E) < p(c_i)$.

Evidence Lift

| # | X | target |
|----|-----|--------|
| 1 | M | c_1 |
| 2 | M | c_1 |
| 3 | M | c_1 |
| 4 | F | c_1 |
| 5 | M | c_0 |
| 6 | F | c_0 |
| 7 | F | c_0 |
| 8 | F | c_0 |
| 9 | M | c_1 |
| 10 | M | c_0 |

| | c_1 | c_0 | |
|-----|-------|-------|---|
| M | 4 | 2 | 6 |
| F | 1 | 3 | 4 |
| | 5 | 5 | |

$$p(c_1) = \frac{1}{2}$$

The prior probability

$$lift_{c_1}(X = F) = \frac{p(X=F|c_1)}{p(X=F)} = \frac{\frac{1}{5}}{\frac{4}{10}} = \frac{1}{2}$$

$$P(c_1|X = F) = p(c_1) \times lift_{c_1}(X = F) = \frac{1}{2} * \frac{1}{2} = \frac{1}{4}$$

Rewritten as $lift_{c_1}(X = F) = \frac{p(c_1|X = F)}{p(c_1)}$

$lift_{c_i}(e_i)$ tells how much more **prevalent** class c_i is in a subpopulation (i.e., $p(c_i|e_i)$), compared to its overall prevalence in general population (i.e., $p(c_i)$).

Evidence Lift of Facebook “Likes”

- ▶ Kosinski et al. (2013) find that what people “like” on FB is predictive of personal traits (usually not apparent):
 - ▶ How they score on psychometric tests (extrovert or conscientious)?
 - ▶ How they score on intelligence tests?
 - ▶ Whether they are (openly) gay?
- ▶ What are the *Likes* that have strong evidence lifts for “high IQ” (IQ>130)?
 - ▶ If a person likes “Lord of the Rings”, the probability he has high-IQ is 69% higher than the baseline (i.e., proportion of high-IQ in general population).

| Like | Lift | Like | Lift |
|------------------------------|------|---------------------------|------|
| <i>Lord Of The Rings</i> | 1.69 | Wikileaks | 1.59 |
| <i>One Manga</i> | 1.57 | Beethoven | 1.52 |
| <i>Science</i> | 1.49 | NPR | 1.48 |
| <i>Psychology</i> | 1.46 | <i>Spirited Away</i> | 1.45 |
| <i>The Big Bang Theory</i> | 1.43 | Running | 1.41 |
| <i>Paulo Coelho</i> | 1.41 | Roger Federer | 1.40 |
| <i>The Daily Show</i> | 1.40 | <i>Star Trek</i> | 1.39 |
| <i>Lost</i> | 1.39 | Philosophy | 1.38 |
| <i>Lie to Me</i> | 1.37 | <i>The Onion</i> | 1.37 |
| <i>How I Met Your Mother</i> | 1.35 | <i>The Colbert Report</i> | 1.35 |
| <i>Doctor Who</i> | 1.34 | <i>Star Trek</i> | 1.32 |
| <i>Howl's Moving Castle</i> | 1.31 | Sheldon Cooper | 1.30 |
| <i>Tron</i> | 1.28 | <i>Fight Club</i> | 1.26 |
| <i>Angry Birds</i> | 1.25 | <i>Inception</i> | 1.25 |
| <i>The Godfather</i> | 1.23 | <i>Weeds</i> | 1.22 |