

## Summer 2021

## 1 Big Idea

## 2 Fundamentals

- Bayes' Theorem
- Bayesian Prediction
- Conditional Independence and Factorization

## 3 Standard Approach: The Naive Bayes' Classifier

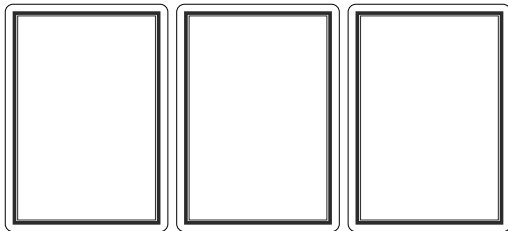
- A Worked Example

## 4 Summary



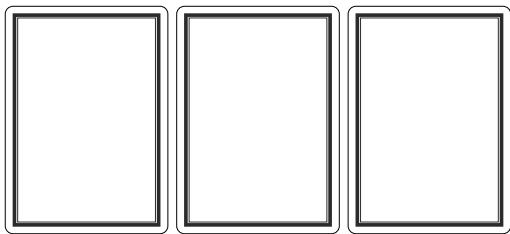


(a)

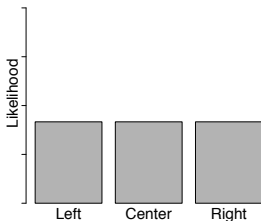


(b)

**Figure:** A game of *find the lady*

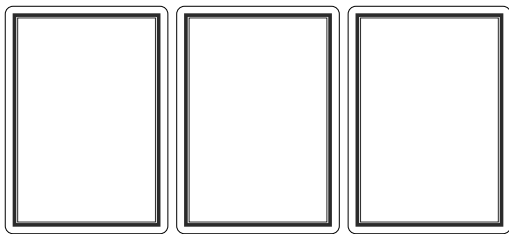


(a)

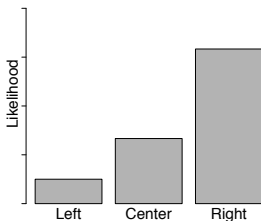


(b)

**Figure:** A game of *find the lady*: (a) the cards dealt face down on a table; and (b) the initial likelihoods of the queen ending up in each position.

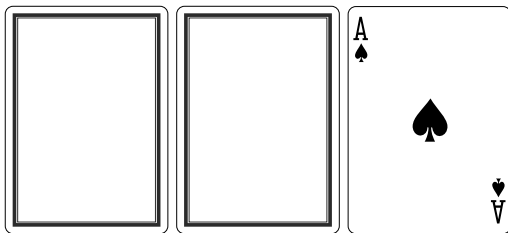


(a)

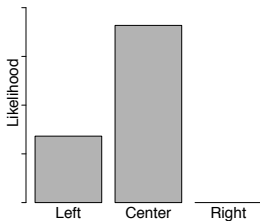


(b)

**Figure:** A game of *find the lady*: (a) the cards dealt face down on a table; and (b) a revised set of likelihoods for the position of the queen based on evidence collected.

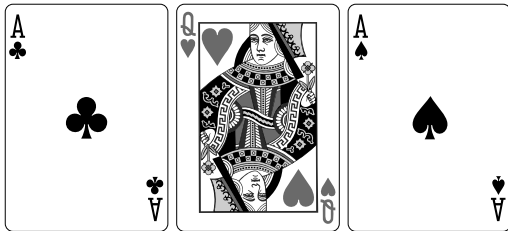


(a)



(b)

**Figure:** A game of *find the lady*: (a) The set of cards after the wind blows over the one on the right; (b) the revised likelihoods for the position of the queen based on this new evidence.



**Figure:** A game of *find the lady*: The final positions of the cards in the game.



## Big Idea

- We can use estimates of likelihoods to determine the most likely prediction that should be made.
- More importantly, we revise these predictions based on data we collect and whenever extra evidence becomes available.

# Fundamentals

**Table:** A simple dataset for MENINGITIS diagnosis with descriptive features that describe the presence or absence of three common symptoms of the disease: HEADACHE, FEVER, and 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

- A **probability function**,  $P()$ , returns the probability of a feature taking a specific value.
- A **joint probability** refers to the probability of an assignment of specific values to multiple different features.
- A **conditional probability** refers to the probability of one feature taking a specific value given that we already know the value of a different feature
- A **probability distribution** is a data structure that describes the probability of each possible value a feature can take. The sum of a probability distribution must equal 1.0.

- A **joint probability distribution** is a probability distribution over more than one feature assignment and is written as a multi-dimensional matrix in which each cell lists the probability of a particular combination of feature values being assigned.
- The sum of all the cells in a joint probability distribution must be 1.0.

$$\mathbf{P}(H, F, V, M) = \begin{bmatrix} P(h, f, v, m), & P(\neg h, f, v, m) \\ P(h, f, v, \neg m), & P(\neg h, f, v, \neg m) \\ P(h, f, \neg v, m), & P(\neg h, f, \neg v, m) \\ P(h, f, \neg v, \neg m), & P(\neg h, f, \neg v, \neg m) \\ P(h, \neg f, v, m), & P(\neg h, \neg f, v, m) \\ P(h, \neg f, v, \neg m), & P(\neg h, \neg f, v, \neg m) \\ P(h, \neg f, \neg v, m), & P(\neg h, \neg f, \neg v, m) \\ P(h, \neg f, \neg v, \neg m), & P(\neg h, \neg f, \neg v, \neg m) \end{bmatrix}$$

- Given a joint probability distribution, we can compute the probability of any event in the domain that it covers by summing over the cells in the distribution where that event is true.
- Calculating probabilities in this way is known as **summing out**.

# Bayes' Theorem

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



100

---

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

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

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

1. 2. 3. 4. 5. 6. 7. 8. 9. 10. 11. 12. 13. 14. 15. 16. 17. 18. 19. 20. 21. 22. 23. 24. 25. 26. 27. 28. 29. 30. 31. 32. 33. 34. 35. 36. 37. 38. 39. 40. 41. 42. 43. 44. 45. 46. 47. 48. 49. 50. 51. 52. 53. 54. 55. 56. 57. 58. 59. 60. 61. 62. 63. 64. 65. 66. 67. 68. 69. 70. 71. 72. 73. 74. 75. 76. 77. 78. 79. 80. 81. 82. 83. 84. 85. 86. 87. 88. 89. 90. 91. 92. 93. 94. 95. 96. 97. 98. 99. 100. 101. 102. 103. 104. 105. 106. 107. 108. 109. 110. 111. 112. 113. 114. 115. 116. 117. 118. 119. 120. 121. 122. 123. 124. 125. 126. 127. 128. 129. 130. 131. 132. 133. 134. 135. 136. 137. 138. 139. 140. 141. 142. 143. 144. 145. 146. 147. 148. 149. 150. 151. 152. 153. 154. 155. 156. 157. 158. 159. 160. 161. 162. 163. 164. 165. 166. 167. 168. 169. 170. 171. 172. 173. 174. 175. 176. 177. 178. 179. 180. 181. 182. 183. 184. 185. 186. 187. 188. 189. 190. 191. 192. 193. 194. 195. 196. 197. 198. 199. 200. 201. 202. 203. 204. 205. 206. 207. 208. 209. 210. 211. 212. 213. 214. 215. 216. 217. 218. 219. 220. 221. 222. 223. 224. 225. 226. 227. 228. 229. 230. 231. 232. 233. 234. 235. 236. 237. 238. 239. 240. 241. 242. 243. 244. 245. 246. 247. 248. 249. 250. 251. 252. 253. 254. 255. 256. 257. 258. 259. 260. 261. 262. 263. 264. 265. 266. 267. 268. 269. 270. 271. 272. 273. 274. 275. 276. 277. 278. 279. 280. 281. 282. 283. 284. 285. 286. 287. 288. 289. 290. 291. 292. 293. 294. 295. 296. 297. 298. 299. 300. 301. 302. 303. 304. 305. 306. 307. 308. 309. 310. 311. 312. 313. 314. 315. 316. 317. 318. 319. 320. 321. 322. 323. 324. 325. 326. 327. 328. 329. 330. 331. 332. 333. 334. 335. 336. 337. 338. 339. 340. 341. 342. 343. 344. 345. 346. 347. 348. 349. 350. 351. 352. 353. 354. 355. 356. 357. 358. 359. 360. 361. 362. 363. 364. 365. 366. 367. 368. 369. 370. 371. 372. 373. 374. 375. 376. 377. 378. 379. 380. 381. 382. 383. 384. 385. 386. 387. 388. 389. 390. 391. 392. 393. 394. 395. 396. 397. 398. 399. 400. 401. 402. 403. 404. 405. 406. 407. 408. 409. 410. 411. 412. 413. 414. 415. 416. 417. 418. 419. 420. 421. 422. 423. 424. 425. 426. 427. 428. 429. 430. 431. 432. 433. 434. 435. 436. 437. 438. 439. 440. 441. 442. 443. 444. 445. 446. 447. 448. 449. 450. 451. 452. 453. 454. 455. 456. 457. 458. 459. 460. 461. 462. 463. 464. 465. 466. 467. 468. 469. 470. 471. 472. 473. 474. 475. 476. 477. 478. 479. 480. 481. 482. 483. 484. 485. 486. 487. 488. 489. 490. 491. 492. 493. 494. 495. 496. 497. 498. 499. 500. 501. 502. 503. 504. 505. 506. 507. 508. 509. 510. 511. 512. 513. 514. 515. 516. 517. 518. 519. 520. 521. 522. 523. 524. 525. 526. 527. 528. 529. 530. 531. 532. 533. 534. 535. 536. 537. 538. 539. 540. 541. 542. 543. 544. 545. 546. 547. 548. 549. 550. 551. 552. 553. 554. 555. 556. 557. 558. 559. 560. 561. 562. 563. 564. 565. 566. 567. 568. 569. 570. 571. 572. 573. 574. 575. 576. 577. 578. 579. 580. 581. 582. 583. 584. 585. 586. 587. 588. 589. 590. 591. 592. 593. 594. 595. 596. 597. 598. 599. 600. 601. 602. 603. 604. 605. 606. 607. 608. 609. 610. 611. 612. 613. 614. 615. 616. 617. 618. 619. 620. 621. 622. 623. 624. 625. 626. 627. 628. 629. 630. 631. 632. 633. 634. 635. 636. 637. 638. 639. 640. 641. 642. 643. 644. 645. 646. 647. 648. 649. 650. 651. 652. 653. 654. 655. 656. 657. 658. 659. 660. 661. 662. 663. 664. 665. 666. 667. 668. 669. 670. 671. 672. 673. 674. 675. 676. 677. 678. 679. 680. 681. 682. 683. 684. 685. 686. 687. 688. 689. 690. 691. 692. 693. 694. 695. 696. 697. 698. 699. 700. 701. 702. 703. 704. 705. 706. 707. 708. 709. 710. 711. 712. 713. 714. 715. 716. 717. 718. 719. 720. 721. 722. 723. 724. 725. 726. 727. 728. 729. 730. 731. 732. 733. 734. 735. 736. 737. 738. 739. 740. 741. 742. 743. 744. 745. 746. 747. 748. 749. 750. 751. 752. 753. 754. 755. 756. 757. 758. 759. 760. 761. 762. 763. 764. 765. 766. 767. 768. 769. 770. 771. 772. 773. 774. 775. 776. 777. 778. 779. 780. 781. 782. 783. 784. 785. 786. 787. 788. 789. 790. 791. 792. 793. 794. 795. 796. 797. 798. 799. 800. 801. 802. 803. 804. 805. 806. 807. 808. 809. 810. 811. 812. 813. 814. 815. 816. 817. 818. 819. 820. 821. 822. 823. 824. 825. 826. 827. 828. 829. 830. 831. 832. 833. 834. 835. 836. 837. 838. 839. 840. 84

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

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

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

- The divisor is the prior probability of the evidence
- This division functions as a normalization constant.

$$0 \leq P(X|Y) \leq 1$$

$$\sum_j P(X_j|Y) = 1.0$$

- 

$$P(Y) = \frac{|\{\text{rows where Y is the case}\}|}{|\{\text{rows in the dataset}\}|}$$

- 

$$P(Y) = \sum_i P(Y|X_i)P(X_i) \quad (1)$$

1. *Journal of Management Studies*, 1997, 34, 1, 1-14.

$$P(t = l | \mathbf{q}[1], \dots, \mathbf{q}[m]) = \frac{P(\mathbf{q}[1], \dots, \mathbf{q}[m] | t = l) P(t = l)}{P(\mathbf{q}[1], \dots, \mathbf{q}[m])}$$

## Chain Rule

$$\begin{aligned}
 P(\mathbf{q}[1], \dots, \mathbf{q}[m]) = \\
 P(\mathbf{q}[1]) \times P(\mathbf{q}[2]|\mathbf{q}[1]) \times \\
 \dots \times P(\mathbf{q}[m]|\mathbf{q}[m-1], \dots, \mathbf{q}[2], \mathbf{q}[1])
 \end{aligned}$$

- To apply the chain rule to a conditional probability we just add the conditioning term to each term in the expression:

$$\begin{aligned}
 P(\mathbf{q}[1], \dots, \mathbf{q}[m] | t = l) = \\
 P(\mathbf{q}[1] | t = l) \times P(\mathbf{q}[2] | \mathbf{q}[1], t = l) \times \dots \\
 \dots \times P(\mathbf{q}[m] | \mathbf{q}[m-1], \dots, \mathbf{q}[3], \mathbf{q}[2], \mathbf{q}[1], t = l)
 \end{aligned}$$

## Bayesian Prediction

| 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     | false | true     | ?          |



$$P(M|h, \neg f, v) = ?$$

- In the terms of Bayes' Theorem this problem can be stated as:

$$P(M|h, \neg f, v) = \frac{P(h, \neg f, v|M) \times P(M)}{P(h, \neg f, v)}$$

- There are two values in the domain of the MENINGITIS feature, 'true' and 'false', so we have to do this calculation twice.

- [illegible]

| 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       |

- $$P(m) = \frac{|\{\mathbf{d}_5, \mathbf{d}_8, \mathbf{d}_{10}\}|}{|\{\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3, \mathbf{d}_4, \mathbf{d}_5, \mathbf{d}_6, \mathbf{d}_7, \mathbf{d}_8, \mathbf{d}_9, \mathbf{d}_{10}\}|} = \frac{3}{10} = 0.3$$
- $$P(h, \neg f, v) = \frac{|\{\mathbf{d}_3, \mathbf{d}_4, \mathbf{d}_6, \mathbf{d}_7, \mathbf{d}_8, \mathbf{d}_{10}\}|}{|\{\mathbf{d}_1, \mathbf{d}_2, \mathbf{d}_3, \mathbf{d}_4, \mathbf{d}_5, \mathbf{d}_6, \mathbf{d}_7, \mathbf{d}_8, \mathbf{d}_9, \mathbf{d}_{10}\}|} = \frac{6}{10} = 0.6$$

- However, as an exercise we will use the chain rule calculate:

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

| 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       |

- Using the chain rule calculate:

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

- 

$$P(m|h, \neg f, v) = \frac{\left( P(h|m) \times P(\neg f|h, m) \right. \\ \left. \times P(v|\neg f, h, m) \times P(m) \right)}{P(h, \neg f, v)} \\ = \frac{0.6666 \times 0.3}{0.6} = 0.3333$$

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

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

- These calculations tell us that it is twice as probable that the patient does not have meningitis than it is that they do even though the patient is suffering from a headache and is vomiting!



## The Paradox of the False Positive

- The mistake of forgetting to factor in the prior gives rise to the **paradox of the false positive** which states that in order to make predictions about a rare event the model has to be as accurate as the prior of the event is rare or there is a significant chance of **false positives** predictions (i.e., predicting the event when it is not the case).

1. *Journal of Management Studies*, 1997, 34, 1, 1-14.

$$\begin{aligned} \mathbb{M}_{MAP}(\mathbf{q}) &= \operatorname{argmax}_{l \in \text{levels}(t)} P(t = l \mid \mathbf{q}[1], \dots, \mathbf{q}[m]) \\ &= \operatorname{argmax}_{l \in \text{levels}(t)} \frac{P(\mathbf{q}[1], \dots, \mathbf{q}[m] \mid t = l) \times P(t = l)}{P(\mathbf{q}[1], \dots, \mathbf{q}[m])} \end{aligned}$$

---

$$\mathbb{M}_{MAP}(\mathbf{q}) = \underset{l \in levels(t)}{\operatorname{argmax}} P(\mathbf{q}[1], \dots, \mathbf{q}[m] \mid t = l) \times P(t = l)$$

| 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       |

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

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

$$P(m \mid h, f, \neg v) = \frac{\left( P(h \mid m) \times P(f \mid h, m) \right. \\ \left. \times P(\neg v \mid f, h, m) \times P(m) \right)}{P(h, f, \neg v)} \\ = \frac{0.6666 \times 0 \times 0 \times 0.3}{0.1} = 0$$

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

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

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

- There is something odd about these results!

## Curse of Dimensionality

As the number of descriptive features grows the number of potential conditioning events grows. Consequently, an exponential increase is required in the size of the dataset as each new descriptive feature is added to ensure that for any conditional probability there are enough instances in the training dataset matching the conditions so that the resulting probability is reasonable.



- The probability of a patient who has a headache and a fever having meningitis should be greater than zero!
- Our dataset is not large enough → our model is **over-fitting** to the training data.
- The concepts of **conditional independence** and **factorization** can help us overcome this flaw of our current approach.

- If knowledge of one event has no effect on the probability of another event, and *vice versa*, then the two events are **independent** of each other.
- If two events  $X$  and  $Y$  are independent then:

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

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

- Recall, that when two event are dependent these rules are:

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

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

- Full independence between events is quite rare.
- A more common phenomenon is that two, or more, events may be independent if we know that a third event has happened.
- This is known as **conditional independence**.

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

X and Y are independent

- If the event  $t = I$  causes the events  $\mathbf{q}[1], \dots, \mathbf{q}[m]$  to happen then the events  $\mathbf{q}[1], \dots, \mathbf{q}[m]$  are conditionally independent of each other given knowledge of  $t = I$  and the chain rule definition can be simplified as follows:

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

- Using this we can simplify the calculations in Bayes' Theorem, under the assumption of conditional independence between the descriptive features given the level  $l$  of the target feature:

$$P(t = l \mid \mathbf{q}[1], \dots, \mathbf{q}[m]) = \frac{\left( \prod_{i=1}^m P(\mathbf{q}[i] \mid t = l) \right) \times P(t = l)}{P(\mathbf{q}[1], \dots, \mathbf{q}[m])}$$

1. *Journal of Management Studies*, 1996, 33, 1, 1-15.

---

1. *Journal of Management Studies*, 1990, 27, 1, 1-14.

$$P(X, Y, Z|W) = \underbrace{P(X|W)}_{Factor1} \times \underbrace{P(Y|W)}_{Factor2} \times \underbrace{P(Z|W)}_{Factor3} \times \underbrace{P(W)}_{Factor4}$$

- The joint probability distribution for the meningitis dataset.

$$\mathbf{P}(H, F, V, M) = \begin{bmatrix} P(h, f, v, m), & P(\neg h, f, v, m) \\ P(h, f, v, \neg m), & P(\neg h, f, v, \neg m) \\ P(h, f, \neg v, m), & P(\neg h, f, \neg v, m) \\ P(h, f, \neg v, \neg m), & P(\neg h, f, \neg v, \neg m) \\ P(h, \neg f, v, m), & P(\neg h, \neg f, v, m) \\ P(h, \neg f, v, \neg m), & P(\neg h, \neg f, v, \neg m) \\ P(h, \neg f, \neg v, m), & P(\neg h, \neg f, \neg v, m) \\ P(h, \neg f, \neg v, \neg m), & P(\neg h, \neg f, \neg v, \neg m) \end{bmatrix}$$



- Assuming the descriptive features are conditionally independent of each other given MENINGITIS we only need to store four factors:

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

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

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

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

$$P(H, F, V, M) = P(M) \times P(H|M) \times P(F|M) \times P(V|M)$$

| 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       |

- Calculate the factors from the data.

$$Factor_1 : < P(M) >$$

$$Factor_2 : < P(h|m), P(h|\neg m) >$$

$$Factor_3 : < P(f|m), P(f|\neg m) >$$

$$Factor_4 : < P(v|m), P(v|\neg m) >$$

$$Factor_4 : \langle P(v|m) = 0.6666, P(v|\neg m) = 0.5714 \rangle$$

$$Factor_1 : \langle P(m) = 0.3 \rangle$$
$$Factor_2 : \langle P(h|m) = 0.6666, P(h|\neg m) = 0.7413 \rangle$$
$$Factor_3 : \langle P(f|m) = 0.3333, P(f|\neg m) = 0.4286 \rangle$$
$$Factor_4 : \langle P(v|m) = 0.6666, P(v|\neg m) = 0.5714 \rangle$$

- Using the factors above calculate the probability of MENINGITIS=*true* for the following query.

| HEADACHE | FEVER | VOMITING | MENINGITIS |
|----------|-------|----------|------------|
| true     | true  | false    | ?          |

$$P(m|h, f, \neg v) = \frac{P(h|m) \times P(f|m) \times P(\neg v|m) \times P(m)}{\sum_i P(h|M_i) \times P(f|M_i) \times P(\neg v|M_i) \times P(M_i)} = \frac{0.6666 \times 0.3333 \times 0.3333 \times 0.3}{(0.6666 \times 0.3333 \times 0.3333 \times 0.3) + (0.7143 \times 0.4286 \times 0.4286 \times 0.7)} = 0.1948$$

$$Factor_1 : \langle P(m) = 0.3 \rangle$$
$$Factor_2 : \langle P(h|m) = 0.6666, P(h|\neg m) = 0.7413 \rangle$$
$$Factor_3 : \langle P(f|m) = 0.3333, P(f|\neg m) = 0.4286 \rangle$$
$$Factor_4 : \langle P(v|m) = 0.6666, P(v|\neg m) = 0.5714 \rangle$$

- Using the factors above calculate the probability of MENINGITIS=*false* for the same query.

| HEADACHE | FEVER | VOMITING | MENINGITIS |
|----------|-------|----------|------------|
| true     | true  | false    | ?          |

$$P(\neg m|h, f, \neg v) = \frac{P(h|\neg m) \times P(f|\neg m) \times P(\neg v|\neg m) \times P(\neg m)}{\sum_i P(h|M_i) \times P(f|M_i) \times P(\neg v|M_i) \times P(M_i)} = \frac{0.7143 \times 0.4286 \times 0.4286 \times 0.7}{(0.6666 \times 0.3333 \times 0.3333 \times 0.3) + (0.7143 \times 0.4286 \times 0.4286 \times 0.7)} = 0.8052$$

$$P(m|h, f, \neg v) = 0.1948$$

$$P(\neg m|h, f, \neg v) = 0.8052$$

- As before, the MAP prediction would be MENINGITIS = 'false'
- The posterior probabilities are not as extreme!



# Standard Approach: The Naive Bayes' Classifier

## Naive Bayes' Classifier

$$\mathbb{M}(\mathbf{q}) = \operatorname{argmax}_{l \in \text{levels}(t)} \left( \prod_{i=1}^m P(\mathbf{q}[i] \mid t = l) \right) \times P(t = l)$$

## Naive Bayes' is simple to train!

- 1 calculate the priors for each of the target levels
- 2 calculate the conditional probabilities for each feature given each target level.

**Table:** A dataset from a loan application fraud detection domain.

| ID | CREDIT HISTORY | GUARANTOR/<br>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 |
| 19 | arrears        | none                      | own          | false |
| 20 | paid           | none                      | own          | false |

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

**Table:** The probabilities needed by a Naive Bayes prediction model calculated from the dataset. Notation key: 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         | ?          |

1

$$P(fr) = 0.3$$

$$P(\neg fr) = 0.7$$

$$P(\text{CH} = \text{'paid'} \mid fr) = 0.1666$$

$$P(\text{CH} = \text{'paid'} \mid \neg fr) = 0.2857$$

$$P(\text{GC} = \text{'none'} \mid fr) = 0.8334$$

$$P(\text{GC} = \text{'none'} \mid \neg fr) = 0.8571$$

$$P(\text{ACC} = \text{'rent'} \mid fr) = 0.3333$$

$$P(\text{ACC} = \text{'rent'} \mid \neg fr) = 0.1429$$

$$\left( \prod_{k=1}^m P(\mathbf{q}[k] \mid fr) \right) \times P(fr) = 0.0139$$

$$\left( \prod_{k=1}^m P(\mathbf{q}[k] \mid \neg fr) \right) \times P(\neg fr) = 0.0245$$

| CREDIT HISTORY | GUARANTOR/COAPPLICANT | ACCOMODATION | FRAUDULENT |
|----------------|-----------------------|--------------|------------|
| paid           | none                  | rent         | ?          |

## A Worked Example

$$P(fr) = 0.3$$

$$P(\neg fr) = 0.7$$

$$P(\text{CH} = \text{'paid'} \mid fr) = 0.1666$$

$$P(\text{CH} = \text{'paid'} \mid \neg fr) = 0.2857$$

$$P(\text{GC} = \text{'none'} \mid fr) = 0.8334$$

$$P(\text{GC} = \text{'none'} \mid \neg fr) = 0.8571$$

$$P(\text{ACC} = \text{'rent'} \mid fr) = 0.3333$$

$$P(\text{ACC} = \text{'rent'} \mid \neg fr) = 0.1429$$

$$\left( \prod_{k=1}^m P(\mathbf{q}[k] \mid fr) \right) \times P(fr) = 0.0139$$

$$\left( \prod_{k=1}^m P(\mathbf{q}[k] \mid \neg fr) \right) \times P(\neg fr) = 0.0245$$

| CREDIT HISTORY | GUARANTOR/COAPPLICANT | ACCOMODATION | FRAUDULENT |
|----------------|-----------------------|--------------|------------|
| paid           | none                  | rent         | 'false'    |



The model is generalizing beyond the dataset!

| ID | CREDIT HISTORY | GUARANTOR/<br>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> |

# Summary

$$P(t|\mathbf{d}) = \frac{P(\mathbf{d}|t) \times P(t)}{P(\mathbf{d})} \quad (2)$$

- A Naive Bayes' classifier naively assumes that each of the descriptive features in a domain is conditionally independent of all of the other descriptive features, given the state of the target feature.
- This assumption, although often wrong, enables the Naive Bayes' model to maximally factorise the representation that it uses of the domain.
- Surprisingly, given the naivety and strength of the assumption it depends upon, a Naive Bayes' model often performs reasonably well.

## 1 Big Idea

## 2 Fundamentals

- Bayes' Theorem
- Bayesian Prediction
- Conditional Independence and Factorization

### 3 Standard Approach: The Naive Bayes' Classifier

- A Worked Example

## 4 Summary