

## Bayes' Law Example

- ▶ **Example:** (1975 British national referendum on whether the UK should remain part of the European Economic Community)
- ▶ Suppose 52% of voters supported the Labour Party and 48% the Conservative Party. Suppose 55% of Labour voters wanted the UK to remain part of the EEC and 85% of Conservative voters wanted this.
- ▶ What is the probability that a person voting “Yes” to remaining in EEC is a Labour voter?

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

## Bayes' Law Example

Note

$$P(Y) = P(Y, L) + P(Y, L^c) = P(Y|L)P(L) + P(Y|L^c)P(L^c).$$

So

$$\begin{aligned} P(L|Y) &= \frac{P(Y|L)P(L)}{P(Y|L)P(L) + P(Y|L^c)P(L^c)} \\ &= \frac{(.55)(.52)}{(.55)(.52) + (.85)(.48)} = 0.41. \end{aligned}$$

## Bayes' Law with Multiple Events

Let  $\mathbf{D}$  represent some observed data and let  $A$ ,  $B$ , and  $C$  be mutually exclusive (and exhaustive) events conditional on  $\mathbf{D}$ . Note that

$$\begin{aligned}P(\mathbf{D}) &= P(A \cap \mathbf{D}) + P(B \cap \mathbf{D}) + P(C \cap \mathbf{D}) \\ &= P(\mathbf{D}|A)P(A) + P(\mathbf{D}|B)P(B) + P(\mathbf{D}|C)P(C).\end{aligned}$$

By Bayes' Law,

$$\begin{aligned}P(A|\mathbf{D}) &= \frac{P(\mathbf{D}|A)P(A)}{P(\mathbf{D})} \\ \Rightarrow P(A|\mathbf{D}) &= \frac{P(\mathbf{D}|A)P(A)}{P(\mathbf{D}|A)P(A) + P(\mathbf{D}|B)P(B) + P(\mathbf{D}|C)P(C)}.\end{aligned}$$

# Bayes' Law with Multiple Events

- ▶ Denoting  $A, B, C$  by  $\theta_1, \theta_2, \theta_3$ , we can write this more generally as

$$P(\theta_i|\mathbf{D}) = \frac{P(\theta_i)P(\mathbf{D}|\theta_i)}{\sum_{j=1}^3 P(\theta_j)P(\mathbf{D}|\theta_j)}.$$

- ▶ If there are  $k$  distinct discrete outcomes  $\theta_1, \dots, \theta_k$ , we have, for any  $i \in \{1, \dots, k\}$ :

$$P(\theta_i|\mathbf{D}) = \frac{P(\theta_i)P(\mathbf{D}|\theta_i)}{\sum_{j=1}^k P(\theta_j)P(\mathbf{D}|\theta_j)},$$

- ▶ The denominator equals  $P(\mathbf{D})$ , the **marginal** distribution of the data.
- ▶ Note if the values of  $\theta$  are portions of the continuous real line, the sum may be replaced by an integral.

# Bayes' Law Example (4 Classes)

**Example:** In the 1996 General Social Survey, for males (age 30+):

- ▶ 11% of those in the lowest income quartile were college graduates.
- ▶ 19% of those in the second-lowest income quartile were college graduates.
- ▶ 31% of those in the third-lowest income quartile were college graduates.
- ▶ 53% of those in the highest income quartile were college graduates.

What is the probability that a college graduate falls in the lowest income quartile?

## Bayes' Law Example (4 Classes)

$$\begin{aligned} P(Q_1|G) &= \frac{P(G|Q_1)P(Q_1)}{\sum_{j=1}^4 P(G|Q_j)P(Q_j)} \\ &= \frac{(.11)(.25)}{(.11)(.25) + (.19)(.25) + (.31)(.25) + (.53)(.25)} = 0.09. \end{aligned}$$

**Exercise:** Find  $P(Q_2|G)$ ,  $P(Q_3|G)$ ,  $P(Q_4|G)$  also. How does this **conditional** distribution differ from the **unconditional** distribution  $\{P(Q_1), P(Q_2), P(Q_3), P(Q_4)\}$ ?

# Statistics Using Bayes' Law

- ▶ We now consider inference about parameters, based on data.
- ▶ Generically denote an unobserved parameter of interest as  $\theta$ .
- ▶ Generically denote our data as  $\mathbf{D}$ .
- ▶ Our probability model for the data, given a value of  $\theta$ , is denoted  $p(\mathbf{D}|\theta)$ .
- ▶ Our model for our prior knowledge about  $\theta$  is denoted  $p(\theta)$ .
- ▶ This could be highly specific or quite vague, depending how uncertain we are about  $\theta$ .

# Statistics Using Bayes' Law

- ▶ We seek to make probability statements about  $\theta$ , **given** some observed data:  $p(\theta|\mathbf{D})$ .
- ▶ By Bayes' Law,

$$p(\theta|\mathbf{D}) = \frac{p(\theta)p(\mathbf{D}|\theta)}{p(\mathbf{D})}.$$

- ▶ Note  $p(\mathbf{D})$  **does not** depend on  $\theta$  and thus carries no information about  $\theta$ .
- ▶ It is simply a **normalizing constant** which makes  $p(\theta|\mathbf{D})$  sum (or integrate) to 1.



# Statistics Using Bayes' Law

- ▶ For inference about  $\theta$ , it is just as good to write

$$p(\theta|\mathbf{D}) \propto p(\theta)p(\mathbf{D}|\theta)$$

- ▶ The LHS is called the **posterior distribution** of  $\theta$  and represents a compromise between the **prior** information about  $\theta$  in  $p(\theta)$  and the information from the sample about  $\theta$  in  $p(\mathbf{D}|\theta)$ .
- ▶ Some useful **summaries** of the posterior are the **posterior mean**

$$E[\theta|\mathbf{D}] = \int \theta p(\theta|\mathbf{D}) d\theta$$

and the **posterior variance**

$$\begin{aligned}\text{var}[\theta|\mathbf{D}] &= E\left\{(\theta - E[\theta|\mathbf{D}])^2|\mathbf{D}\right\} \\ &= \int (\theta - E[\theta|\mathbf{D}])^2 p(\theta|\mathbf{D}) d\theta \\ &= \int \theta^2 p(\theta|\mathbf{D}) d\theta - 2E[\theta|\mathbf{D}] \int \theta p(\theta|\mathbf{D}) d\theta \\ &\quad + \left(E[\theta|\mathbf{D}]\right)^2 \int p(\theta|\mathbf{D}) d\theta \\ &= E[\theta^2|\mathbf{D}] - \left(E[\theta|\mathbf{D}]\right)^2\end{aligned}$$

- ▶ If the values of  $\theta$  are discrete, sums would replace the integrals.

CHAPTER 2 SLIDES START HERE

## Some Notation

- ▶ **Notation:** We hereby denote our data as the  $n \times k$  matrix  $\mathbf{X}$ .
- ▶ We denote the parameter(s) of interest (possibly multidimensional) to be the vector  $\theta$ .
- ▶ We will denote our posterior distribution for  $\theta$  using  $\pi(\cdot)$ .

# Likelihood Theory

- ▶ The likelihood function  $L(\theta|\mathbf{X})$  is a function of  $\theta$  that shows how “likely” are various parameter values  $\theta$  to have produced the data  $\mathbf{X}$  that **were observed**.
- ▶ In classical statistics, the specific value of  $\theta$  that maximizes  $L(\theta|\mathbf{X})$  is the maximum likelihood estimator (MLE) of  $\theta$ .
- ▶ In many common probability models, when the sample size  $n$  is large,  $L(\theta|\mathbf{X})$  is unimodal in  $\theta$ .
- ▶ **Note:** Unlike  $p(\theta|\mathbf{X})$ ,  $L(\theta|\mathbf{X})$  does **not necessarily** obey the usual laws for probability distributions.
- ▶ Also, in the classical framework, all the randomness within  $L(\theta|\mathbf{X})$  is attached to  $\mathbf{X}$ , not to  $\theta$ .

# Likelihood Theory

- ▶ Mathematically, if the data  $\mathbf{X}$  represent iid observations from probability distribution  $p(\mathbf{X}|\theta)$ , then

$$L(\theta|\mathbf{X}) = \prod_{i=1}^n p(\mathbf{X}_i|\theta)$$

(where  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are the  $n$  data vectors).

- ▶ The **Likelihood Principle** of Birnbaum states that (given the data) all of the evidence about  $\theta$  is contained in the likelihood function.
- ▶ Likelihood Principle implies: Two experiments that yield equal (or proportional) likelihoods should produce equivalent inference about  $\theta$ .

# The Bayesian Framework

- ▶ Suppose we observe an iid sample of data  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$ .
- ▶ Now  $\mathbf{X}$  is considered fixed and known.
- ▶ We also **must** specify  $p(\boldsymbol{\theta})$ , the prior distribution for  $\boldsymbol{\theta}$ , based on any knowledge we have about  $\boldsymbol{\theta}$  **before** observing the data.
- ▶ Our model for the distribution of the data will give us the likelihood

$$L(\boldsymbol{\theta}|\mathbf{X}) = \prod_{i=1}^n p(\mathbf{X}_i|\boldsymbol{\theta}).$$

# The Bayesian Framework

- ▶ Then by Bayes' Law, our posterior distribution is

$$\begin{aligned}\pi(\boldsymbol{\theta}|\mathbf{X}) &= \frac{p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{X})}{p(\mathbf{X})} \\ &= \frac{p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{X})}{\int_{\Theta} p(\boldsymbol{\theta})L(\boldsymbol{\theta}|\mathbf{X}) d\boldsymbol{\theta}}\end{aligned}$$

- ▶ Note that the **marginal distribution** of  $\mathbf{X}$ ,  $p(\mathbf{X})$ , is simply the joint density  $p(\boldsymbol{\theta}, \mathbf{X})$  (i.e., the numerator) with  $\boldsymbol{\theta}$  integrated out.
- ▶ With respect to  $\boldsymbol{\theta}$ , it is simply a **normalizing constant** that ensures that  $\pi(\boldsymbol{\theta}|\mathbf{X})$  integrates to 1.



# The Bayesian Framework

- ▶ Since  $p(\mathbf{X})$  carries no information about  $\theta$ , for conciseness we may drop it and write

$$\pi(\theta|\mathbf{X}) \propto p(\theta)L(\theta|\mathbf{X}).$$

- ▶ Often we can calculate the posterior distribution by multiplying the prior by the likelihood and **then** normalizing the posterior at the **last** step, by including the necessary constant.

# Bayesian Inference: Posterior Intervals

- ▶ Simple values like the posterior mean  $E[\boldsymbol{\theta}|\mathbf{X}]$  and posterior variance  $var[\boldsymbol{\theta}|\mathbf{X}]$  can be useful in learning about  $\boldsymbol{\theta}$ .
- ▶ Quantiles of  $\pi(\boldsymbol{\theta}|\mathbf{X})$  (especially the posterior median) can also be a useful summary of  $\boldsymbol{\theta}$ .
- ▶ The ideal summary of  $\boldsymbol{\theta}$  is an interval (or region) with a certain probability of containing  $\boldsymbol{\theta}$ .
- ▶ Note that a classical (frequentist) **confidence interval** does not exactly have this interpretation.

# Definitions of Coverage

- ▶ **Defn.:** A random interval  $(L(\mathbf{X}), U(\mathbf{X}))$  has  $100(1 - \alpha)\%$  **frequentist coverage** for  $\theta$  if, **before** the data are gathered,

$$P[L(\mathbf{X}) < \theta < U(\mathbf{X})|\theta] = 1 - \alpha.$$

(**Pre-experimental**  $1 - \alpha$  coverage)

- ▶ Note that if we observe  $\mathbf{X} = \mathbf{x}$  and plug  $\mathbf{x}$  into our confidence interval formula,

$$P[L(\mathbf{x}) < \theta < U(\mathbf{x})|\theta] = \begin{cases} 0 & \text{if } \theta \notin (L(\mathbf{x}), U(\mathbf{x})) \\ 1 & \text{if } \theta \in (L(\mathbf{x}), U(\mathbf{x})) \end{cases}$$

(**NOT** Post-experimental  $1 - \alpha$  coverage)

# Definitions of Coverage

- ▶ **Defn.:** An interval  $(L(\mathbf{x}), U(\mathbf{x}))$ , based on the observed data  $\mathbf{X} = \mathbf{x}$ , has  $100(1 - \alpha)\%$  **Bayesian coverage** for  $\theta$  if

$$P[L(\mathbf{x}) < \theta < U(\mathbf{x}) | \mathbf{X} = \mathbf{x}] = 1 - \alpha.$$

(**Post-experimental**  $1 - \alpha$  coverage)

- ▶ The frequentist interpretation is less desirable if we are performing inference about  $\theta$  based on a **single** interval.