

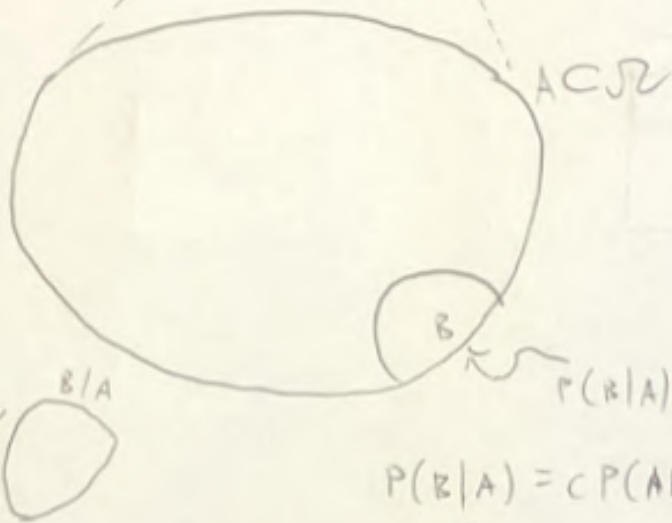
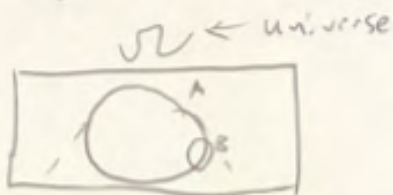
Lecture 4 - 2/10/21

Let A : smoking and B : lung cancer

$$P(A) = 0.200$$

$$P(B) = 0.060$$

$$P(A|B) = 0.036$$


$$P(\text{lung cancer} | \text{Smoking})$$

$$P(B|A)$$

given, conditional universe
(B, relative to A)

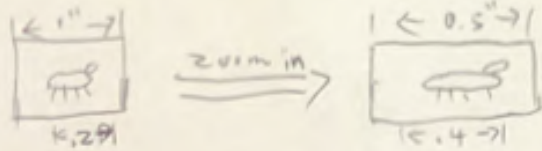
"proportional"

$$P(B|A) \propto P(AB)$$

$c = \text{scale factor}$

$$P(B|A) = \frac{P(AB)}{P(A)} = \frac{P(\sqrt{2})}{P(A)} P(AB) = \frac{P(AB)}{P(A)}$$

Definition of conditional probability

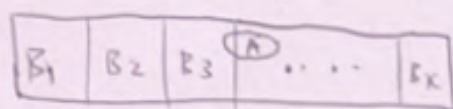


$$P(A|B) = \frac{P(A \cap B)}{P(B)} \Rightarrow P(A \cap B) = P(A|B) P(B) \Rightarrow P(B|A) = \frac{P(A \cap B) P(B)}{P(A)}$$

Bayes rule

$$P(A) = P(AB) + P(AB^c) = P(A|B)P(B) + P(A|B^c)P(B^c)$$

If B_1, B_2, \dots, B_K are mutually exclusive and collectively exhaustive
 $\forall i \neq j \quad B_i \cap B_j = \emptyset$ $\mathcal{B} = \bigcup_{i=1}^K B_i$

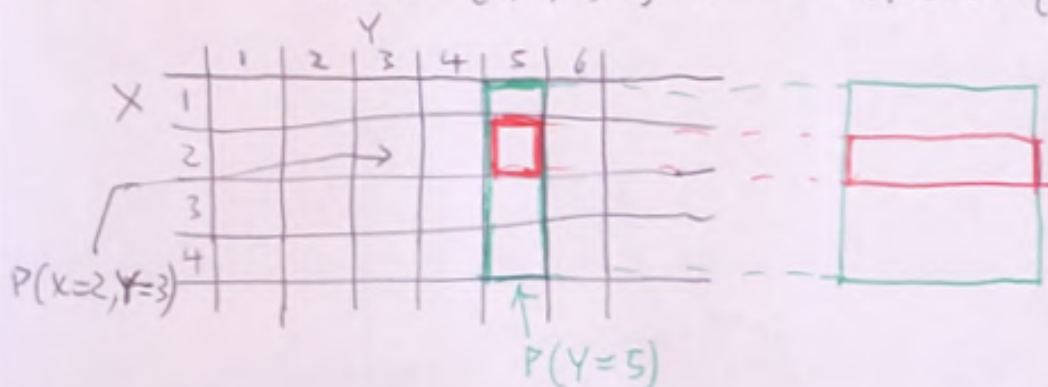


$$P(A) = \sum_{k=1}^K P(A, B_k) = \sum_{k=1}^K P(A|B_k)P(B_k)$$

"margining out the B_k 's or integrating out the B_k 's"

$$P(B_i|A) = \frac{P(A, B_i)}{P(A)} = \frac{P(A|B_i)P(B_i)}{\sum_{k=1}^K P(A|B_k)P(B_k)} \quad \text{Bayes Theorem}$$

Bayes Rule and Bayes Thm for r.v.'s. Imagine two r.v.'s X, Y and the $\text{Supp}[X] = \{1, 2, 3, 4\}$ and $\text{Supp}[Y] = \{1, 2, 3, 4, 5, 6\}$



$$P(Y=5) = P(Y=5, X=1) + P(Y=5, X=2) + P(Y=5, X=3) + P(Y=5, X=4)$$

$$= \sum_{X \in \text{Supp}[X]} P(Y=5, X=x)$$

$$P(X=2|Y=5) = \frac{P(X=2, Y=5)}{P(Y=5)}$$

$$P(Y) := P(Y=y) = \sum_{X \in \text{Supp}[X]} P(Y=y, X=x) = \sum_x p(y, x)$$

$$P(X=x|Y=y) = P(X=x, Y=y) = \frac{p(x, y)}{P(y)} \quad \leftarrow \text{JMF}$$

Conditional PMF

$$f(y) = \int_{\text{Supp}[X]} f(x, y) dx$$

$$f_{X|Y}(x, y) = \frac{f(x, y)}{f(y)} \quad \leftarrow \text{JDF}$$

Back to the story... can we use Bayes Rule to tell us anything about inference for parameter θ given data x ($x = \langle x_1, \dots, x_n \rangle$).

Consider:
$$P(\theta|x) = \frac{P(x|\theta)P(\theta)}{P(x)}$$

$P(x_1, \dots, x_n|\theta) = \mathcal{L}(\theta; x_1, \dots, x_n) \rightarrow$ JMF, but called "likelihood"

What is wrong with this equation? Previously, θ , the unknown parameter was assumed to be a fixed real value. Thus, $\theta \sim \text{Deg}(\theta)$. Then, this equation is trivial. If you plug in the actual value of $\theta = \theta_0$ on the right hand side then you get:

$$P(\theta = \theta_0|x) = \frac{P(x|\theta)(1)}{\sum_{\theta \in \Theta} P(x|\theta)P(\theta)} = \frac{P(x|\theta_0)}{P(x|\theta_0)} = 1$$

$$P(\theta \neq \theta_0|x) = \frac{\cancel{P(x|\theta_0)}(0)}{\sum_{\theta \in \Theta} P(x|\theta)P(\theta)} = \frac{0}{P(x|\theta_0)} = 0$$

This was a mean exam problem but not super interesting since you don't know θ_0 and even if you did, this doesn't help with the three goals of inference.

The big leap: let θ be a r.v.! Then $P(\theta)$ has a distribution (either discrete or continuous). But it's a constant! This is the big philosophical problem in Bayesian Statistics / Bayesian Inference. Some authors say it's still a constant but $P(\theta)$ represents uncertainty in its value. Purists say that's nonsense.

$$P(\theta|x) = \frac{\overbrace{P(x|\theta)}^{\text{likelihood}} \overbrace{P(\theta)}^{\text{prior}}}{\underbrace{P(x)}_{\text{prior predictive distribution}}}$$

$P(x) \begin{cases} \rightarrow = \sum_{\theta \in \Theta} P(x|\theta)P(\theta) \text{ if } \theta \text{ discrete} \\ \rightarrow = \int_{\Theta} P(x|\theta)P(\theta)d\theta \text{ if } \theta \text{ contin.} \end{cases}$

prior: thoughts summed up in a distribution over Θ the parameter space **before** seeing any data. There is no x within it.

Frequentists say this is "subjective" and not real!

posterior: thoughts summed up in a distribution over Θ the parameter space **after** seeing the data x which is why it's conditional on x !

Notation for the rest of class: "p" now discrete PMF/conditional mass function **or** continuous PDF/conditional density function. I won't use "f" anymore.

$X = \text{iid Bernoulli}, x = \{0, 1\}$ $P(X|\theta) = \theta^x(1-\theta)^{1-x}$

Let $\Theta = \{0.5, 0.75\} \neq (0, 1)$

$$P(\theta = 0.75|x) \stackrel{?}{>} P(\theta = 0.5|x)$$

$$P(\theta = 0.75|x) \stackrel{\text{Bayes Thm}}{=} \frac{P(x|\theta = 0.75) P(\theta = 0.75)}{P(x|\theta = 0.5) P(\theta = 0.5) + P(x|\theta = 0.75) P(\theta = 0.75)}$$

$$P(x|\theta = 0.75) = (0.75)^2(0.25) = 0.141, \quad P(x|\theta = 0.5) = (0.5)^3 = 0.125$$

We need $P(\theta = 0.75)$ and $P(\theta = 0.5)$ to complete the calculations. That's the prior, $P(\theta)$. It's subjective. What do you think it should be? Amir says $P(\theta = 0.75) = 0.2$ and $P(\theta = 0.5) = 0.8$ because he feels that way.

An automatic rule is called the "principle of indifference" (Laplace's idea so it's sometimes called the "Laplace prior"). This principle says that all values of θ in the parameter space are equally likely. In our case,

$$P(\theta) = \begin{cases} 0.5 & \text{if } \theta = 0.7 \\ 0.5 & \text{if } \theta = 0.5 \end{cases} \quad \text{In general, } P(\theta) = \frac{1}{|\Theta|}, \text{ this formula only works for finite parameter spaces.}$$

$$P(\theta = 0.75|x) = \frac{(0.141)(0.5)}{(0.125)(0.5) + (0.141)(0.5)} = 0.53^*$$

* $\theta = 0.75$ is your point estimate.

$$P(\theta = 0.5|x) = \frac{(0.125)(0.5)}{(0.125)(0.5) + (0.141)(0.5)} = 0.47$$

$$P(\theta = 0.75) = 0.5 \quad \xrightarrow{x} \quad P(\theta = 0.75|x) = 0.53$$

This is called Bayesian Conditionalism.