The Generalization Jigsaw

BAYESIANISM

- (methods)
- Parameter confidence
- Model choice: model weighting

FREQUENTISM

- (methods)
- Parameter confidence
- Model choice: hypothesis testing

EMPIRICISM

- Evaluating model fit
- Model choice: holdout evaluation

PUTTING THE JIGSAW TOGETHER

Bayes's rule, done right

For two discrete random variables X and Y,

$$\mathbb{P}(X = x | Y = y) = \frac{\mathbb{P}(X = x)\mathbb{P}(Y = y | X = x)}{\mathbb{P}(Y = y)} \quad \text{when } \mathbb{P}(Y = y) > 0$$

For two discrete or continuous random variables X and Y,

$$\Pr_{X}(x|Y=y) = \frac{\Pr_{X}(x) \Pr_{Y}(y|X=x)}{\Pr_{Y}(y)} \quad \text{when } \Pr_{Y}(Y) > 0$$

What do these "conditional likelihoods" even mean?

Joint distribution

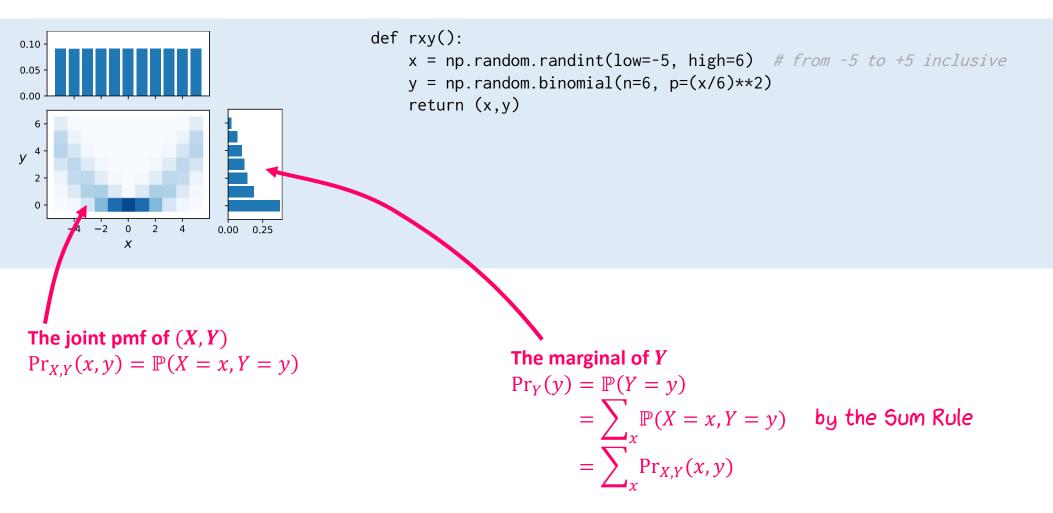
```
def rxy():
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     y = np.random.binomial(n=6, p=(x/6)**2)
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        return (x,y)
by definite conditional probability

The joint pmf of (X,Y)

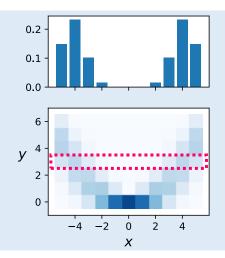
P(X=x) P(Y=y | X=x) = \frac{1}{1!} {\binom{6}{y}} {\binom{x}{6}}^y {\binom{x}{6}}^y {\binom{x}{6}}^y

P(X=x) P(Y=y | X=x) = \frac{1}{1!} {\binom{6}{y}} {\binom{x}{6}}^y {\binom{x}{6}}^y
```

Marginal random variables



Conditional random variables



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

QUESTION. What is X conditional on Y = 3?

$$\mathbb{P}(X = x | Y = 3) = \frac{\mathbb{P}(X = x, Y = 3)}{\mathbb{P}(Y = 3)} = \frac{\Pr_{X,Y}(x,3)}{\Pr_{Y}(3)}$$

i.e. take the Y=3 row, then rescale it to sum to 1 We can think of "X conditional on Y = 3" as a random variable ...

We've provided a valid probability mass function:

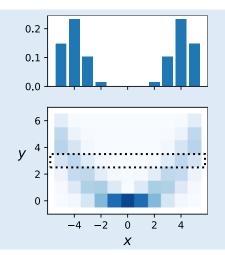
$$pmf_{s}(\cdot)>0$$
 $\sum_{x}pmf_{s}(x)=1$

 $pmf_3(.) > 0$ $\sum_{x} pmf_3(x) = 1$ Sample space: $\Omega = \{-5, -4, \dots, 4, 5\}$

Code to generate values from it:

$$\begin{array}{lll} \text{def rx_given_y():} & \text{def rx_given_y():} \\ & \text{while True:} & \Omega = \{-5, \ldots, 5\} \\ & \text{x,y = rxy()} & \text{p = [pmf}_3(\text{x}) \text{ for x in } \Omega] \\ & \text{if y == 3: break} & \text{return np.random.choice}(\Omega, \text{p=p}) \\ & \text{return x} \end{array}$$

Conditional random variables



```
def rxy():
    x = np.random.randint(low=-5, high=6) # from -5 to +5 inclusive
    y = np.random.binomial(n=6, p=(x/6)**2)
    return (x,y)
```

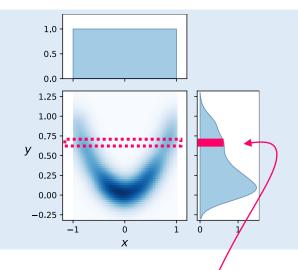
We *define* the **conditional random** variable, written (X|Y=y), by specifying its likelihood:

$$Pr_{(X|Y=y)}(x) = \frac{Pr_{X,Y}(x,y)}{Pr_{Y}(y)}$$

This likelihood is also written $Pr_X(x|Y=y)$.

i.e. take the Y=y row from the joint pmf, then rescale it

Conditional random variables (continuous case)



```
def rxy():
    x = np.random.uniform(-1,1)
    y = np.random.normal(loc=x**2, scale=0.1)
    return (x,y)
```

QUESTION. What is X conditional on Y = 0.6?

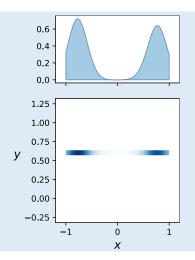
Take the Y=0.6 slice of the joint pdf, then rescale it so it integrates to I i.e. so we get a legitimate pdf.

(Rescale it by dividing by $Pr_{Y}(0.6)$) where $Pr_{Y}(\cdot)$ is the marginal for Y.)

The marginal for *Y*

$$\Pr_{Y}(y) = \int_{x} \Pr_{X,Y}(x,y) \ dx$$

Conditional random variables (continuous case)



```
def rxy():
    x = np.random.uniform(-1,1)
    y = np.random.normal(loc=x**2, scale=0.1)
    return (x,y)
```

QUESTION. What is X conditional on Y = 0.6?

Take the Y=0.6 slice of the joint pdf, then rescale it so it integrates to I i.e. so we get a legitimate pdf.

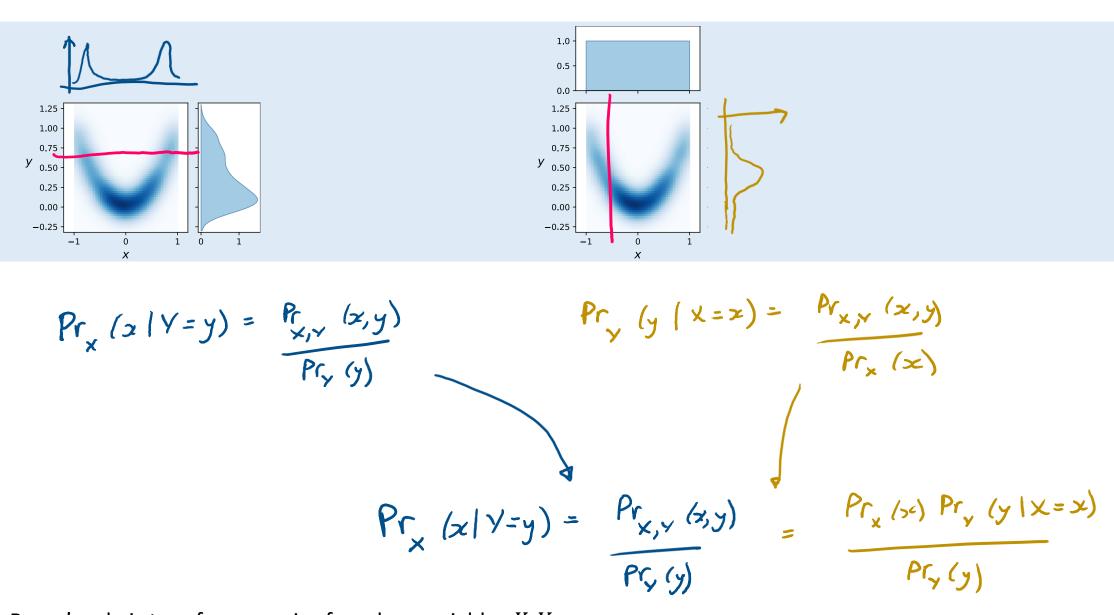
(Rescale it by dividing by $Pr_{Y}(0.6)$) where $Pr_{Y}(\cdot)$ is the marginal for Y.)

We define the conditional random variable

$$(X|Y = y)$$
 by specifying its likelihood:

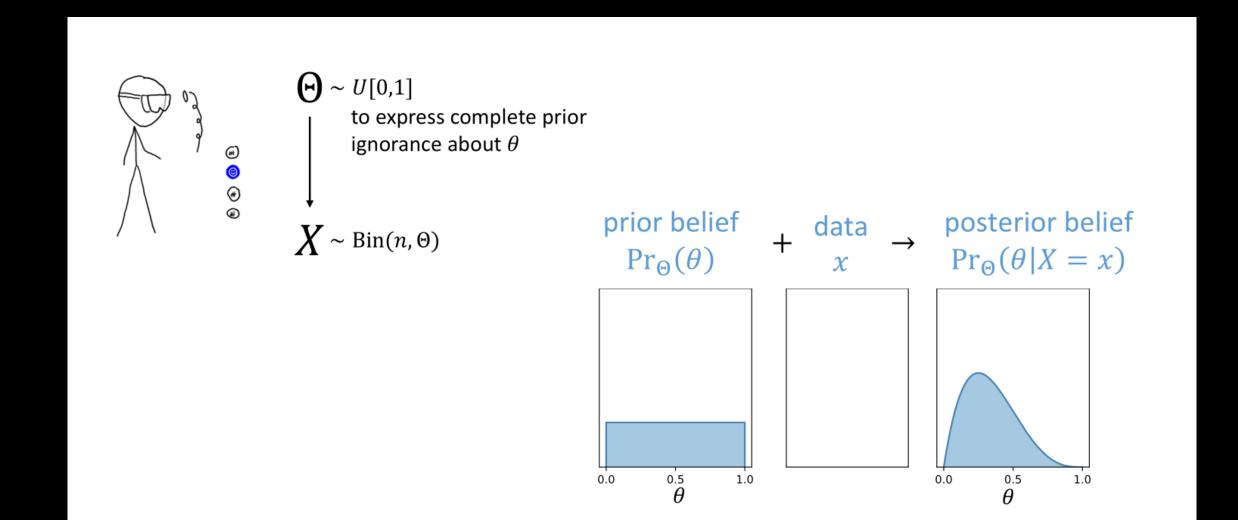
$$Pr_X(x|Y = y) = \frac{Pr_{X,Y}(x,y)}{Pr_Y(y)}$$

Bayes's rule



Bayes's rule is true for any pair of random variables X, Y. It's only useful for "sequential models" i.e. when the question tells us $Pr_X(x)$ and $Pr_Y(y|X=x)$.

Now we have the tech to apply Bayes's rule to problems with continuous random variables.



Bayes's rule for discrete or continuous random variables

For two random variables X and Y,

$$Pr_X(x|Y=y) = \frac{Pr_X(x) Pr_Y(y|X=x)}{Pr_Y(y)}$$
 when $Pr_Y(y) > 0$

In practice, we write it as

$$\Pr_{X}(x|Y=y) = \kappa \Pr_{X}(x) \Pr_{Y}(y|X=x)$$

$$\Pr_{X}(x|Y=y)$$

then figure out κ so that $\Pr_X(\cdot | Y = y)$ is a legitimate likelihood function

Choose
$$k$$
 so that
$$\int_{z} k \, Pr_{x}(x) \, Pr_{y}(y|X=x) \, dx = 1$$

Exercise.

Consider the pair of random variables (Θ, X) where

$$\Theta \sim U[0,1], \qquad X \sim Bin(4,\Theta)$$

Find the distribution of
$$(\Theta|X=1)$$
.

$$\Pr_{\Theta}(\theta) = 1$$

$$\Pr_{X}(x|\Theta=\theta) = \binom{1}{x} \Theta^{x} (I-\Theta)^{x-x} = 4 \Theta (I-\Theta)^{x} (x - 1)^{x-x}$$

$$\Pr_{X}(x|\Theta=\theta) = \kappa \Pr_{\Theta}(\theta) \Pr_{X}(1|\Theta=\theta) = \kappa \times 1 \times 4 \Theta (I-\Theta)^{x}$$

$$= \kappa' \Theta (I-\Theta)^{x} \times 4 \otimes (I-\Theta)^{$$

Exercise.

Consider the pair of random variables (Θ, X) where

$$\Theta \sim U[0,1], \qquad X \sim Bin(4,\Theta)$$

Find the distribution of $(\Theta|X=1)$.

$$Pr_{\Theta}(\theta|X=1) = \kappa Pr_{\Theta}(\theta) Pr_X(1|\Theta=\theta)$$

$$= K \Theta (1-\Theta)^{3} = K B(\alpha,\beta)$$

$$= K B(\alpha,\beta)$$
So this constant

Beta

Probability density function	
Notation	Beta (α, β)
PDF	$x^{\alpha-1}(1-x)^{\beta-1}$
	$\overline{B(\alpha,\beta)}$
1 1 16	where $B(\alpha,\beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ and Γ
ndard pdf	is the Gamma function.
	Notation

$$\frac{\Theta^{\alpha-1}(1-0)^{\beta-1}}{B(\alpha,\beta)}$$

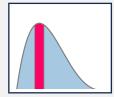
where
$$\alpha = 2$$
, $\beta = 4$

must be 1 (otherwise this pdf wouldn't integrate to 1 wr.t. θ)

Thus
$$(\Theta(X=1) \sim Beta(X=2, \beta=4)$$

It's easier to communicate the posterior distribution by reporting a confidence interval.

E.g. what is
$$\mathbb{P}(\Theta \in [.2, .3] \mid X = 1)$$
?

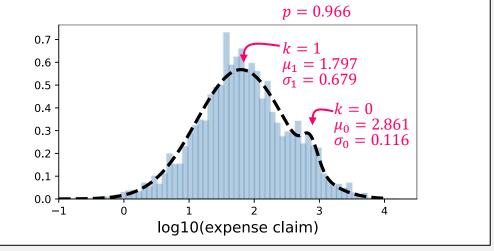


Exercise 5.2.3 (Bayesian classification)

In a dataset of MP expense claims, let y_i be \log_{10} of the claim amount in record i. A histogram of the y_i suggests we use a Gaussian mixture model with two components,

$$K = \begin{cases} 1 & \text{with prob } p \\ 0 & \text{with prob } 1 - p \end{cases}$$
$$Y \sim \text{Normal}(\mu_K, \sigma_K^2)$$

Find the probability that a claim amount £5000 belongs to the component k=0.



In other words, find
$$\mathbb{P}(K=0|Y=y)$$
, where $y=\log_{10}5000=3.70$.

$$Pr_K(k) =$$

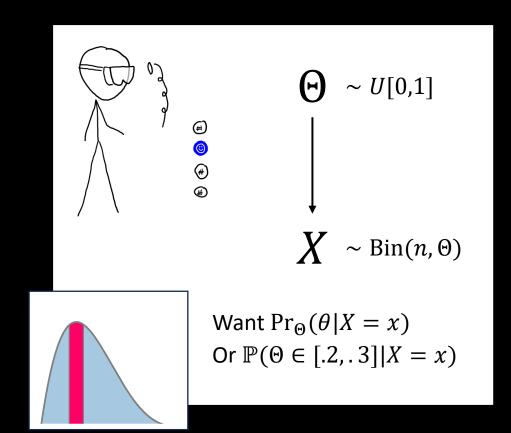
$$\Pr_{Y}(y|K=k) =$$

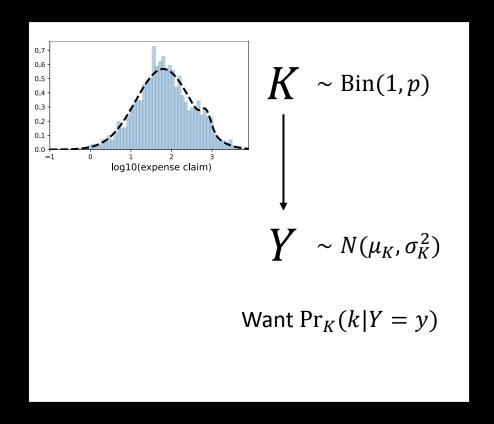
$$\Pr_K(k|Y=y) = \kappa \Pr_K(k) \Pr_Y(y|K=k)$$

$$\Pr_{K}(k) \Pr_{Y}(y|K=k)$$

$$\text{Choose } K \quad \text{s.t.} \quad \sum_{k} k \; \Pr_{K}(k) \; \Pr_{S}(y|k=k) = 1$$

We've used Bayes's rule to find the posterior distribution of random variable. We've had tricky integrals / sums, to find the normalizing constant.





If we want to give our answer in a more practical form, as a confidence interval, there's another tricky integral.

Next: how we can use sampling and computation instead of tricky integrals

§6. Computational methods

What's the chance that a randomly thrown dart will hit the mystery object A?



Let X be the location of a randomly thrown dart, and let x_1, \dots, x_n be some throws.

The probability of hitting A is

The probability of hitting
$$A$$
 is
$$\mathbb{P}(X \in A) \approx \frac{1}{n} \sum_{i=1}^{n} 1_{x_i \in A}$$

```
1 # Let X \sim N(\mu = 1, \sigma = 3). What is \mathbb{P}(X > 5)?
2 x = np.random.normal(loc=1, scale=3, size=10000)
3 i = (x > 5) — i is • vec of 10000 sould
4 np.mean(i) what fraction of them one True?
```

Expectation

For a real-valued random variable *X*

$$\mathbb{E}X = \begin{cases} \sum_{x} x \Pr_{X}(x), & \text{if } X \text{ is discrete} \\ \int_{x} x \Pr_{X}(x) dx, & \text{if } X \text{ is continuous} \end{cases}$$

Law of the Unconscious Statistician

For a random variable X and a real-valued function h

$$\mathbb{E}h(X) = \begin{cases} \sum_{x} h(x) \operatorname{Pr}_{X}(x), & \text{if } X \text{ is discrete} \\ \int_{x} h(x) \operatorname{Pr}_{X}(x) dx, & \text{if } X \text{ is continuous} \end{cases}$$

If we want to know the average properties of a rich random variable (e.g. the average length of a random tweet), we have to use real-valued property readout functions h(X) so that we can take averages.

Monte Carlo integration

$$\mathbb{E}h(X) \approx \frac{1}{n} \sum_{i=1}^{n} h(x_i)$$

where x_1, \dots, x_n is a sample drawn from X

Let
$$h(x) = 1 \times eA$$

By Monthe Conform

$$Eh(x) \approx \int_{1}^{\infty} \int_{1}^{\infty} h(x_{i})$$

Let $y = h(x)$

Let $y = h(x)$

$$EY = 0 \times P(y = 0) + 1 \times P(y = 1)$$

$$= P(Y = 1)$$

$$= P(1 \times eA = 1)$$

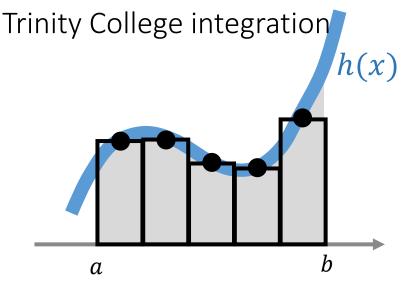
$$= P(X \in A)$$

Let X be the location of a randomly thrown dart, and let x_1, \dots, x_n be some throws.

The probability of hitting A is

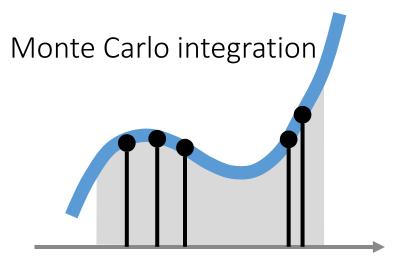
$$\mathbb{P}(X \in A) \approx \frac{1}{n} \sum_{i=1}^{n} 1_{x_i \in A}$$





$$\int_{x=a}^{b} h(x) dx \approx \sum_{i=1}^{n} h(x_i) \frac{b-a}{n}$$

where x_i is the midpoint of interval i



Let's instead approximate this integral using Monte Carlo. Let $X \sim U[a,b]$. By Monte Carlo,

$$\mathbb{E}h(X) \approx \frac{1}{n} \sum_{i=1}^{n} h(x_i) \text{ where } x_1, \dots, x_n \text{ sampled from } X$$

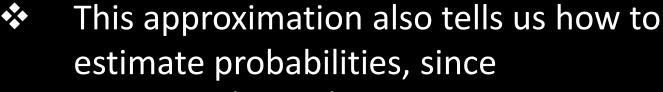
$$\int_{x=a}^{b} h(x) \Pr_X(x) dx = \int_{x=a}^{b} h(x) \frac{1}{b-a} dx$$

Thus,
$$\int_{x=a}^{b} h(x) dx \approx \frac{b-a}{n} \sum_{i=1}^{n} h(x_i)$$

COMPUTATIONAL METHODS

If we want $\mathbb{E}h(X)$ but the maths is too complicated, we can approximate $\mathbb{E}h(X) \approx n^{-1} \sum_{i=1}^{n} h(X_i)$

$$\mathbb{E}h(X) \approx n^{-1} \sum_{i=1}^{n} h(x_i)$$
 where x_1, \dots, x_n are sampled from X



$$\mathbb{P}(X \in A) = \mathbb{E}1_{X \in A}$$

In Bayesian analysis, our aim is to estimate probabilities, e.g. $\mathbb{P}(\Theta \in [.2, .3] | X = 1)$. How can we do this computationally?

