

# Lecture 7 probabilistic models - Bayesian methods

wbg231

December 2022

## 1 introduction

- we have so far done frequentest probabilistic models using MLE
- we are going to use Bayesian methods to get some uncertainty around the prediction

## 2 classical stats

### 2.1 parametric family of densities

- a parametric family of densities is a set

$$\{p(y|\theta) : \theta \in \Theta\}$$

this is a set of distributions

- where  $P(y|\theta)$  is a density on a sample space  $y$ , and  $\theta$  is a parameter in a parameter space  $\Theta$
- this is a common starting point for Bayesian statistics.

### frequentest statistics

- in classical of frequentest statistics we are working with a parametric family of distributions

$$\{P(y|\theta : \theta \in \Theta)\}$$

- however we assume that there is some  $\theta_{true} \in \Theta$  which has governed the distribution of our observed data
- so if we know  $\theta_{true}$  there would be no need for statistics
- but we can not view the true data generating process, as we only have a finite sample  $\mathcal{D} : (y_1 \dots y_n)$  generated independent and identically distributed from  $P(y|\theta_{true})$

## point estimation

- one type of statistical problem is point estimation
- a **statistic**  $s = s(\mathcal{D})$  is any function of our data
- a **point estimator of  $\theta$**  is a function of our data  $\hat{\theta} = \hat{\theta}(\mathcal{D})$  for some  $\theta \in \Theta$
- a good point estimate will have  $\hat{\theta} \approx \theta_{true}$
- we want a point estimate to have the following properties
  1. **consistency** that is if we think of our point estimate on  $n$  data points as  $\hat{\theta}_n$  then we have  $\lim_{n \rightarrow \infty} \hat{\theta}_n \rightarrow \theta_{true}$  that is as we get more data our point estimate generally more accurate, that is as we get more data our error gets lower
  2. **efficiency** formally the efficiency of an unbiased estimator  $\hat{\theta}$  of parameter  $\theta$  is  $e(t) = \frac{1}{\frac{I(\theta)}{var(\hat{\theta})}}$ , where  $I(\theta)$  is a measure of the amount about how much our observed data  $\mathcal{D} = \{y_1 \dots y_n\}$  tells us about our parameter  $\theta$  it related to entropy which is how predicible a variable is effectively. basically this is saying we want our estimator to train well on relatively little data efficiency
- maximum likelihood estimators as consistent and efficient under reasonable assumptions

## 2.2 coin example

- we have a parametric family of mass functions  $P(\text{heads}|\theta) = \theta$  for some  $\theta \in \Theta := (0, 1)$

## 2.3 coin flipping mle

- suppose we have dataset  $\mathcal{D} = \{H \dots T\}$  where we have  $n_h$  heads and  $n_t$  tails and the flips are iid
- the likelihood function of our data is then  $\mathcal{L}_{\mathcal{D}}(\theta) = P(\mathcal{D}|\theta) = P(y_1 \dots y_n|\theta) = \prod_{i=1}^n P(y_i|\theta) = (\theta)^{n_h} (1 - \theta)^{n_t}$
- we can max the log likelihood of our data as a function of  $\theta$  as  $\max_{\theta \in \Theta} \ell(\theta) = \max_{\theta \in \Theta} n_h \log(\theta) + n_t \log(1 - \theta)$
- we can then see that  $\frac{\partial \ell}{\partial \theta} = \frac{n_h}{\theta} - \frac{n_t}{1 - \theta}$
- setting this equal to zero we get  $\theta_{mle} = \frac{n_h}{n_h + n_t}$  which is the empirical fraction of heads which makes sense

## bayesian statistics

- in bayesian stats we introduce a prior distribution
- **the prior distribution** is defined as  $P(\theta)$ , and it represents our beliefs about how  $\theta$  is distributed over the parameter space  $\Theta$  prior to seeing any data

### a bayesian model

- there are two pieces to a bayesian model
  1. a parametric family of densities  $P(\mathcal{D}|\theta \in \Theta)$  that is basically a set of distributions we think our data given the parameter may have
  2. we also need our prior  $P(\theta) : \theta \in \Theta$
- given both of these pieces we can write our joint density  $P(\mathcal{D}, \theta) = P(\mathcal{D}|\theta)P(\theta)$  so we have the joint density of our data and the model

- **the posterior**

- **the posterior distribution** for  $\theta$  is  $P(\mathcal{D}|\theta)$
- the prior is our beliefs about the parameter before seeing any data
- the posterior represents how we rationally update our beliefs about  $\theta$  after seeing our data
- we can write the posterior as  $P(\mathcal{D}|\theta) = \frac{P(\mathcal{D}, \theta)}{P(\theta)} = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})}$  where we think of both sides as a function of  $\theta$  for a fixed dataset  $\mathcal{D}$
- given our data set is fixed  $P(\mathcal{D})$  is constant so we can ignore it
- so in practice we solve for  $P(\mathcal{D}|\theta)P(\theta)$  as we know that  $P(\theta|\mathcal{D}) \propto P(\mathcal{D}, \theta)P(\theta)$

### coin flipping bayesian example

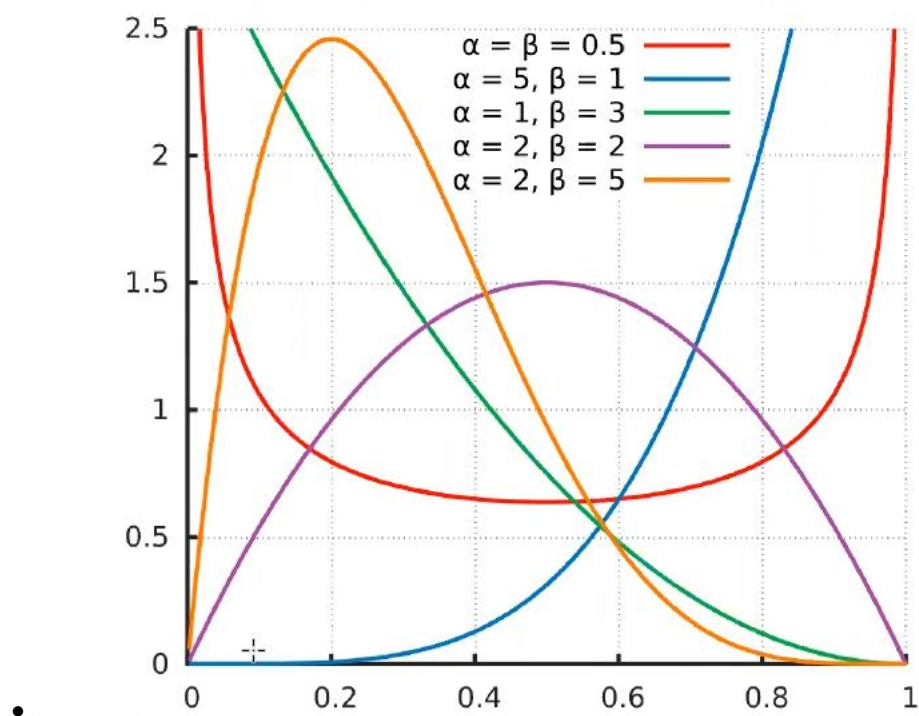
- suppose we have a parametric family of mass functions  $P(\text{heads}|\theta) = \theta$  where  $\theta \in \Theta = (0, 1)$
- we need a prior distribution  $P(\theta)$  on  $\Theta$
- typically we choose a distribution from the beta family

## beta distributions

- given we assume our prior is  $\theta \sim \text{beta}(\alpha, \beta)$  we know that  $P(\theta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$
- the beta family takes two parameters
- given these parameters we get a distribution over  $\theta$  but notice that this distribution is independent of our data
- the shape of the beta distribution can vary a lot.

$$\theta \sim \text{Beta}(\alpha, \beta)$$

$$p(\theta) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1}$$



- the beta distribution is nice because it is only defined in 0,1
- given that  $\theta \sim \text{beta}(\alpha, \beta)$
- $E[\theta] = \frac{\alpha}{\alpha+\beta}$  the proof is kind of tedious so i am just going to link it

- then we can find the mode of the beta distribution. the mode is the most common value in a pdf that is  $\text{argmax}_{\theta \in [0,1]} P(\theta) = \frac{\alpha-1}{\alpha+\beta-2}$
- showing the mode is a big less tedious classically just note  $P(\theta) = \frac{1}{B(a,b)} \theta^a (1-\theta)^b$ , then we can take the derivative and solve for  $\theta^*$

### coin flipping prior

- lets say that our prior is  $\theta \sim \text{beta}(c, c)$  where  $h = c = t$  that is we are assuming that the likelihood of heads and tails are equal then our mean and mode are both equal to  $\frac{1}{2}$

### coin flipping posterior

- so our likelihood of the data given  $\theta$  is  $\mathcal{L}_{\mathcal{D}}(\theta) = P(\mathcal{D}|\theta) = \theta^{n_h} (1-\theta)^{n_t}$
- then our posterior density is

$$P(\theta|\mathcal{D}) \propto P(\mathcal{D}|\theta)P(\theta) \propto \theta^{n_h} (1-\theta)^{n_t} \times (1-\theta)^t \theta^h = \theta^{h-1+n_h} (1-\theta)^{t-1+n_t}$$

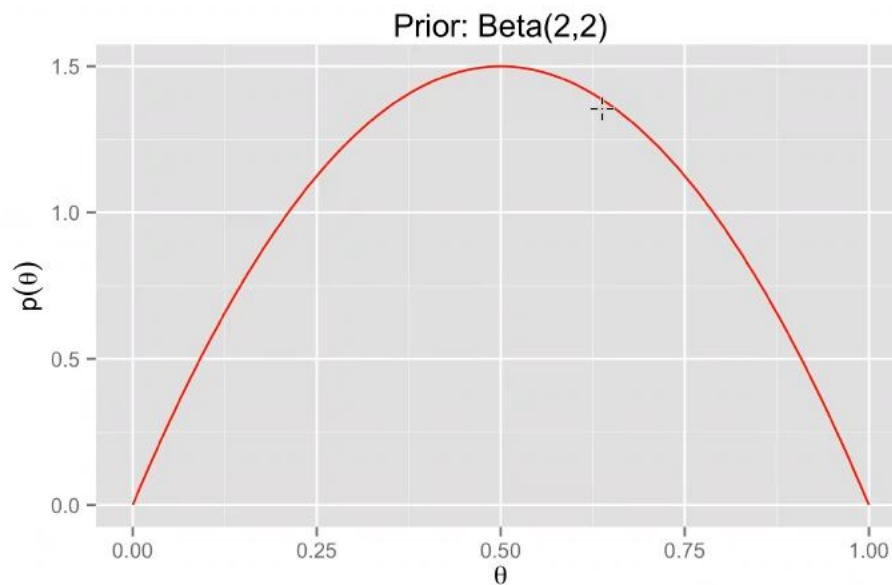
- note that our posterior is in the beta family that is  $P(\theta|d) \propto \theta^{h-1+n_h} (1-\theta)^{t-1+n_t} \Rightarrow \theta|\mathcal{D} \sim \text{beta}(h+n_h, t+n_t)$
- as the number of coin flips goes to infinity the prior will matter less as the values of h and t are fixed and  $n_h, n_t$  grow, so when we have a lot of data we weigh it quite heavily
- and when we do not have much data we rely more on our prior

### conjugate priors

- in this case the posterior was in the same family of distributions as the prior
- this makes the math easy
- let  $\pi$  be a family of posterior distributions on  $\Theta$
- let  $P$  a parametric family of distributions on parameter space  $\Theta$
- family of distributions  $\pi$  is **conjugate** to the parametric model P if  $\forall p(\theta) \in \pi$  (that is all priors on  $\pi$ ) the posterior is in  $\pi$  that is  $P(\theta) \in \pi \Rightarrow P(\theta|\mathcal{D}) \in \pi$
- the beta family is conjugate to a bernoulli model
- this is not easy to do, but it is nice when it works

### coin flipping concrete example

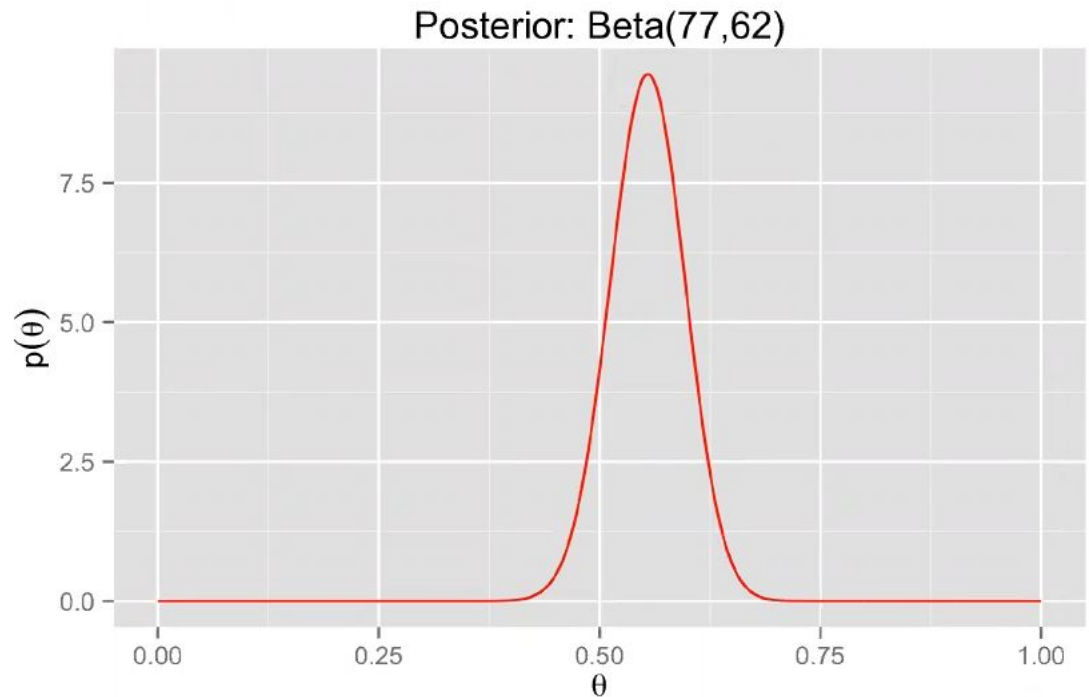
- suppose we have a parametric probabilistic model of our coin  $P(\text{heads}|\theta) = \theta$
- and our parameter space is  $\theta \in \Theta = [0, 1]$
- and our prior is  $\theta \sim \text{beta}(2, 2)$



- 
- our prior assumes that  $\theta$  will be distributed kind of evenly but centered at 0

### with data

- suppose we have some data, where we saw 75 heads and 60 tails
- doing maximum likelihood estimation we would get  $\hat{\theta}_{mle} = .556$
- using bayesian methods we would get a posterior  $\theta|\mathcal{D} \sim \text{beta}(77, 62)$



- 
- as we can see both estimates are centered at about 55
- but doing it with bayesian methods we have a distribution of parameters as opposed to just a single estimate of  $\theta$

### bayesian point estimates

- what if want to give a point estimate from our posterior
- there are a few common options
  1. posterior mean  $\hat{\theta} = E[\theta|\mathcal{D}]$
  2. maximum a posteriori estimate (MAP)  $\hat{\theta} = \text{argmax}_{\theta} P(\theta|\mathcal{D})$  (this is the mode of the posterior )

### what else can we do with a prior

- we can use it to quantify our uncertainty around the estimate
- we could make a credible set for  $\theta$  that is an interval  $[a, b] : P(\theta \in [a, b]|\mathcal{D}) \geq \alpha$  this is effectively bayesian confidence interval

- we could also select a point estimate using bayesian decisions theory. this requires us to chose a loss function, and chose an action to minimize the expected risk with respect to

### 3 bayesian decision theory

- we need the following ingredients
  1. parameter space  $\Theta$
  2. prior  $p(\theta) : \theta \in \Theta$
  3. action space  $A$
  4. loss function  $\ell : A \times \theta \Rightarrow \mathbb{R}$
- the output is no longer the parameter it is an action
- the posterior risk of an action  $a \in A$  (also alled tthe expected loss under the posterior )is

$$r(a) =: E[\ell(\theta, a)|\mathcal{D}] = \int \ell(\theta, a)P(\theta|\mathcal{D})d\theta$$

- so that is more or less looking at a weighted average of our action over all possible values of theta
- this is more robust as we are only choosing an action a, not a point estimate of our parameter
- a byes action  $a^* = \min_{a \in A} r(a)$  the action that minimizes risk, ie the best action

#### bayesian point estimation

- we have a data set  $\mathcal{D}$  genrated by  $P(y|\theta)$  for some unkown  $\theta \in \Theta$
- we want a point estiamte for  $\theta$
- we need to chose a prior  $P(\theta) : \theta \in \Theta$  and loss  $\ell(\hat{\theta}, \theta)$
- and our goal is to fund the action  $\hat{\theta}$  that minimizes the posterior risk.
- the point of this is that bayesian point estimation can be looked with the bayesian decision theory framework



## important cases

- squared loss  $\ell(\theta - \hat{\theta}) = (\theta - \hat{\theta})^2$  minimizing this gives the  $\hat{\theta}$  = posterior mean
- zero one loss  $\ell(\theta - \hat{\theta}) = \mathbf{1}(\theta \neq \hat{\theta})$  minimizing this gives the  $\hat{\theta}$  = posterior mode (not a good idea for when  $\theta$  is continuous)
- $\ell(\theta - \hat{\theta}) = |\theta - \hat{\theta}|$  minimizing this gives the  $\hat{\theta}$  = posterior median

## squared loss example

- want to find an action  $\hat{\theta} \in \Theta : \hat{\theta} \in \operatorname{argmin}_{\theta \in \Theta} \int (\theta - \hat{\theta})^2 P(\theta|\mathcal{D}) d\theta$
- we can take the partial of our risk as  $\frac{\partial r(\hat{\theta})}{\partial \hat{\theta}} = -2 \int (\theta - \hat{\theta}) p(\theta|\mathcal{D}) d\theta = -2 \int \theta P(\theta|\mathcal{D}) d\theta + 2\hat{\theta} \int P(\theta|\mathcal{D}) d\theta = -2 \int \theta P(\theta|\mathcal{D}) d\theta + 2\hat{\theta}$
- setting that equal to zero yields  $\hat{\theta} = \int \theta P(\theta|\mathcal{D}) d\theta = E[\theta|\mathcal{D}]$
- so in other words the optimal action according to squared loss is the posterior mean
- all inferences and actions can be taken with only a prior and loss function
- in the bayesian approach we do not need to justify our estimator, we just need to specify our family of distributions and the prior
- try to use the conjugate prior when you can
- for a lot of data the prior does not matter very much

## recap of conditional probabilistic models

### conditional probabilistic models

- input space  $X$
- outcome space  $y$
- action space  $A = \{P(y)|p \text{ is a probability distribution on } y\}$
- hypothesis space  $\mathcal{F}$  contains prediction function mapping  $f : X \rightarrow A$
- prediction function  $f \in \mathcal{F} : f(x)$  produces a distribution on  $y$
- a parametric family of conditional densities is a set  $\{P(y|x, \theta) : \theta \in \Theta\}$
- where  $P(y|x, \theta)$  is a density on the outcome space  $y$  for each  $x$  in input space and  $\theta$  is a parameter in the parameter space
- the action space here is a probability distribution not just a decision

## likelihood function

- we can as always find our likelihood function given our data set  $P(\mathcal{D}|x_1 \dots x_n, \theta) = \prod_{i=1}^n P(y_i|x_i, \theta)$
- the mle estimator is the the  $\hat{\theta}_{mle} = \operatorname{argmax}_{\theta \in \Theta} \mathcal{L}_{\mathcal{D}}(\theta)$
- the corresponding prediction function is  $f(x) = P(y|x, \hat{\theta}_{mle})$

## bayesian conditiional models

- input space  $X = \mathbb{R}^D$  outcome space  $Y = \mathbb{R}$
- parametric family of distributions  $\{P(y|x, \theta) : \theta \in \Theta\}$
- prior  $P(\theta) : \theta \in \Theta$  (so we have added a prior )
- the posterior is  $P(\theta|\mathcal{D}, X) \propto P(\mathcal{D}|\theta)P(\theta) = \mathcal{L}_{\mathcal{D}}(\theta)p(\theta)$
- we don't worry about the denominator in this case because the dataset is fixed so we can just look at the proportional
- we can use bayesian decisions theory to derive point estimates, we have a few choices for how we do this. also depends on our loss function

## bayesian prediction function

- we want to find a prediction function that takes input  $x \in X$  and produces a distribution on  $Y$
- in the frequentest approach we chose a conditional family of probability distributions (hypothesis space) and select one conditional probability from the family based on some rule like the mle
- in Bayesian setting we chose a parametric family of conditional densities

$$\{P(y|x, \theta) : \theta \in \Theta\}$$

and a prior distribution  $P(\theta) : \theta \in \Theta$

- with a bayesian model how do we predict a distribution on  $y$  for input  $x$
- we do not need to make a discrete selection from the hypothesis space: we can maintain some uncertainty

### **prior predictive distribution**

- suppose we have not yet observed any data
- in the Bayesian setting we can still produce a prediction function
- call **the prior predictive distribution**

$$x \rightarrow P(y|x) = \int P(y|x, \theta) p(\theta) d\theta$$

- this is an average of all conditional densities in our family weighted by the prior.
- so we are considering all possible  $\theta$  and we get out a  $P(y|x)$

### **posterior predictive distribution**

- after seeing our data set  $\mathcal{D}$
- the posterior predictive distribution is given by

$$x \rightarrow P(y|x, \mathcal{D}) = \int P(y|x, \theta) P(\theta|\mathcal{D}) d\theta$$

- this is an average of all conditional densities in our hypothesis space weighted by the posterior distribution, so we are again considering all  $\theta$
- we have not chosen a particular  $\theta$  we consider all and just weighted by their likelihood

### **comparing to the frequentist approach**

- in Bayesian stats we have two distributions on  $\Theta$  the prior, and the posterior  $P(\theta|\mathcal{D})$
- these distributions are over parameters corresponding to the distribution on the hypothesis space so, what we get out is a distribution on the hypothesis space
- in the frequentist approach we just pick one  $\hat{\theta}$  that we think is best
- in the Bayesian approach we weight over all possible outcomes

### **what if we don't want a full distribution on y**

- once we have a predictive distribution  $P(y|x, \mathcal{D})$  we can generate a single prediction
- there are many choices depending on what loss we want to minimize

## gaussian regression example

### example in 1 dimension

- pick up at 47 minutes