# Applied Machine Learning

Maximum Likelihood and Bayesian Reasoning
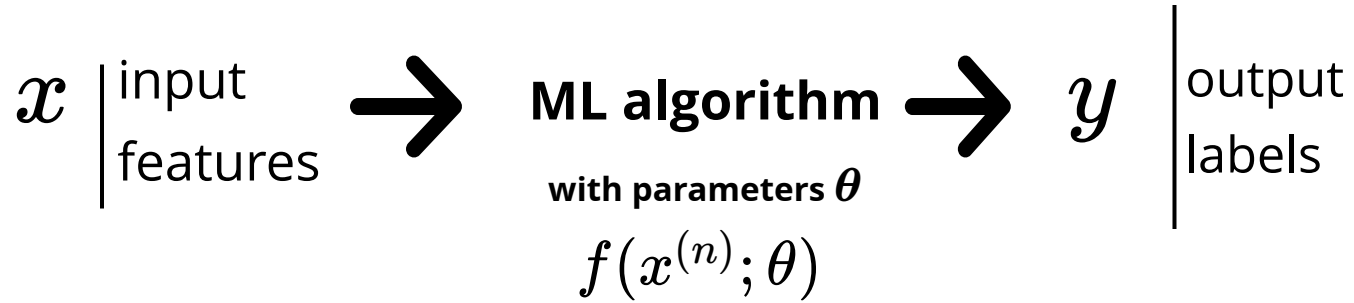
**Oumar Kaba**

# Admin

- Add/drop deadline is tomorrow

- Do the quizz before the tomorrow if you are unsure about your math background

- We will solve the issue with study groups later this week

- Office hours for this week: after each class

- Bonus points for lecture notes/summaries

# Model fitting

$x$ | input features  $\longrightarrow$  **ML algorithm**  $\longrightarrow$  $y$ | output labels

**with parameters $\theta$**

$$f(x^{(n)}; \theta)$$

The process of estimating the model parameters $\theta$ from given data $\mathcal{D}$, is the core of training ML models which often boils down into optimization of an loss function $\mathcal{L}(\theta)$
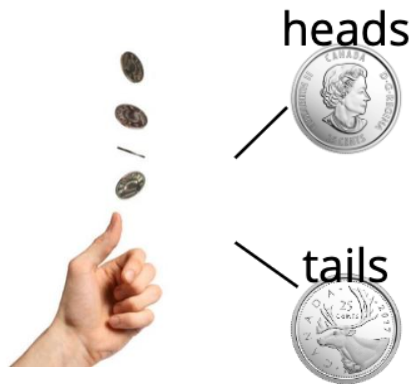
Where do these loss functions come from?

Often from **maximum likelihood** or **Bayesian** methods

# Case study

Fundamental machine learning problem that we will study today

You are given observations (e.g. data) of coin flips from a possibly rigged coin



$$\mathcal{D} = \{0, 0, 0, 0, 1\}$$

**Coin flip is just one example, could be anything binary:**

- Someone purchasing product or not
- Someone getting infected by covid or not
- Bus arrives on time or not
- A penalty kick is scored or not
- Social media post is liked or not
- etc.

What is your **estimate** for the probability of the next throw being head (1) or tail (0)?

# Objectives

learn common parameter estimation methods and understand what it means to learn a probabilistic model of the data

- using maximum likelihood principle
- using Bayesian inference
  - prior, posterior, posterior predictive
  - MAP inference
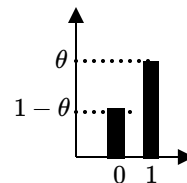  - Beta-Bernoulli conjugate pairs

# Parameter estimation

we suppose a coin's head/tail outcome has a **Bernoulli distribtion**

$$\mathrm{Bernoulli}(x|\theta) = \theta^x (1-\theta)^{(1-x)}$$

reminder: Bernoulli random variable takes values of 0 or 1, e.g. head/tail in a coin toss

$$p(x|\theta) = \begin{cases} \theta & x = 1 \\ 1 - \theta & x = 0 \end{cases}$$

this is our **probabilistic model** of some head/tail IID data   $\mathcal{D} = \{0, 0, 1, 1, 0, 0, 1, 0, 0, 1\}$

**Objective:** learn the model parameter $\theta$

if we are only interested in the counts, we can also use **Binomial distribution**

$$\mathrm{Binomial}(N, N_h|\theta) = \binom{N}{N_h} \theta^{N_h} (1-\theta)^{N-N_h}$$

$|\mathcal{D}|$

\# heads $N_h = \sum_{x \in \mathcal{D}} x$

$N_t$

6

# Maximum likelihood

a coin's head/tail outcome has a **Bernoulli distribtion**

$$\text{Bernoulli}(x|\theta) = \theta^x(1-\theta)^{(1-x)}$$

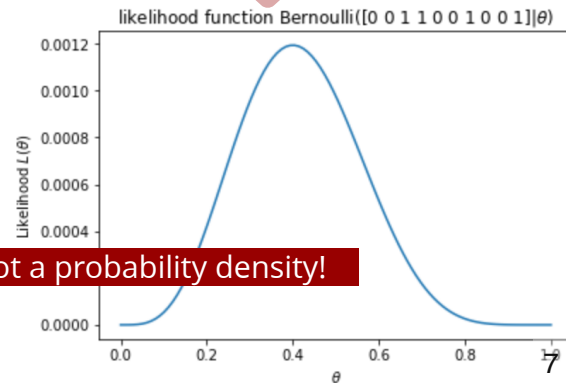this is our **probabilistic model** of some head/tail IID data $\mathcal{D} = \{0, 0, 1, 1, 0, 0, 1, 0, 0, 1\}$

**Objective: learn** the model parameter $\theta$

**Idea:** find the parameter $\theta$ that maximizes the probability of observing $\mathcal{D}$

**Max-likelihood** assignment

**Likelihood**  $L(\theta; \mathcal{D}) = \prod_{x \in \mathcal{D}} \text{Bernoulli}(x|\theta) = \theta^4(1-\theta)^6$  is a function of $\theta$

*pick the parameters that assign the highest probability to the training data*


likelihood function Bernoulli([0 0 1 1 0 0 1 0 0 1]|θ)

note that this is not a probability density!

# Maximizing log-likelihood

**likelihood** $L(\theta; \mathcal{D}) = \prod_{x \in \mathcal{D}} p(x; \theta)$

*using product here creates extreme values*

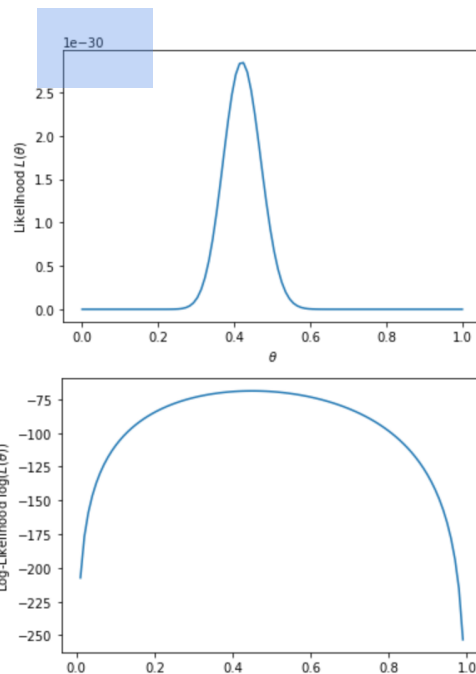*for 100 samples in our example, the likelihood shrinks below 1e-30*

**log-likelihood** has the same maximum but it is well-behaved

$$\ell(\theta; \mathcal{D}) = \log(L(\theta; \mathcal{D})) = \sum_{x \in \mathcal{D}} \log(p(x; \theta))$$

how do we find the max-likelihood parameter? $\theta^* = \arg\max_\theta \ell(\theta; \mathcal{D})$

*for some simple models we can get the **closed form solution***

*for complex models we need to use **numerical optimization***



8

# Maximizing log-likelihood



**log-likelihood** $\ell(\theta; \mathcal{D}) = \log(L(\theta; \mathcal{D})) = \sum_{x \in \mathcal{D}} \log(\text{Bernoulli}(x; \theta))$

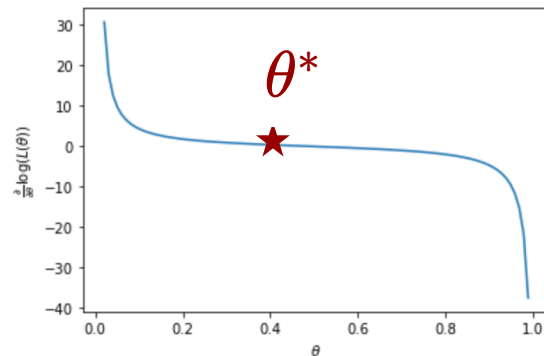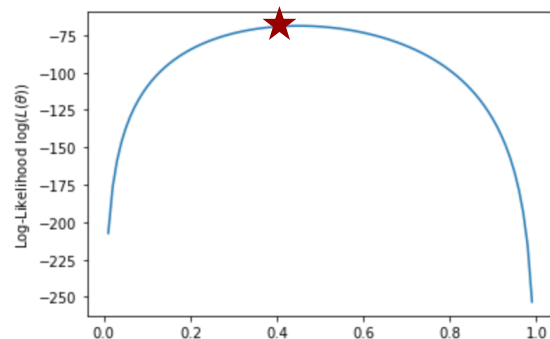**observation:** at maximum, the derivative of $\ell(\theta; \mathcal{D})$ is zero

**idea:** set the the derivative to zero and solve for $\theta$

> **example**   max-likelihood for Bernoulli

$$\frac{\partial}{\partial \theta} \ell(\theta; \mathcal{D}) = \frac{\partial}{\partial \theta} \sum_{x \in \mathcal{D}} \log\left(\theta^x (1-\theta)^{(1-x)}\right)$$
$$= \frac{\partial}{\partial \theta} \sum_x x \log \theta + (1-x) \log(1-\theta)$$
$$= \sum_x \frac{x}{\theta} - \frac{1-x}{1-\theta} = 0$$



which gives $\theta^{MLE} = \frac{\sum_{x \in \mathcal{D}} x}{|\mathcal{D}|}$ is simply the portion of heads in our dataset

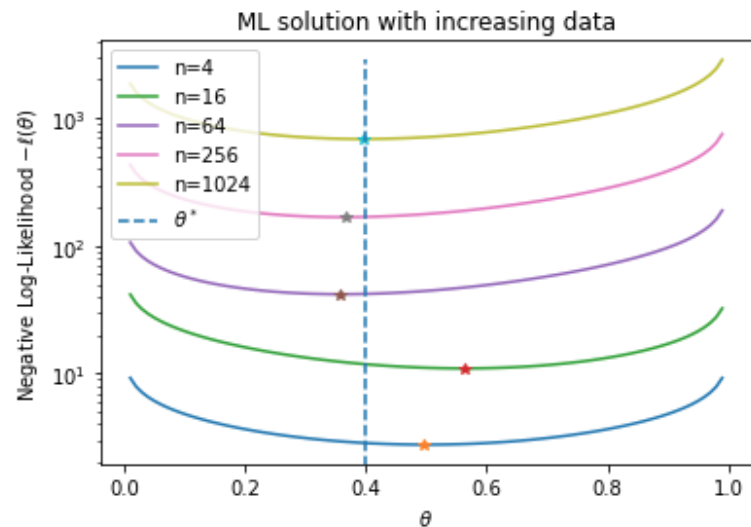what is $\theta^{MLE}$ when $\mathcal{D} = \{0, 0, 1, 1, 0, 0, 1, 0, 0, 1\}$?

# Problem with maximum likelihood

max-likelihood estimate does not reflect our uncertainty:

- e.g. for $\mathcal{D} = \{1\}$, $\theta^{MLE} = 1$. If we observe only one head, predicts all future tosses are head!
- e.g., $\theta^{MLE} = .2$ for both 1/5 heads and 1000/5000 heads

    - in which case are we more certain of the predicted $\theta$?

    How can we quantify our uncertainty
    about our prediction?



ML solution with increasing data

# Bayesian approach



How can we quantify our uncertainty about our prediction?
capture it using a conditional probability distribution instead of a single best guess

Using the Bayesian inference approach

- we maintain a *distribution* over parameters $p(\theta)$    prior    what do we believe about $\theta$ before any observation
- after observing $\mathcal{D}$ we update this distribution $p(\theta|\mathcal{D})$    posterior

how to update degree of certainty given data?   using **Bayes rule**

$$\underset{\text{observed}}{\overset{\text{hidden}}{p(\theta|\mathcal{D})}} = \frac{\overset{\text{prior}}{p(\theta)}p(\mathcal{D}|\theta)}{p(\mathcal{D})}$$

**likelihood** of the data
*previously denoted by* $L(\theta; \mathcal{D})$

**evidence**: this is a normalization, **marginal likelihood of data**

$$p(\mathcal{D}) = \int p(\theta')p(\mathcal{D}|\theta')\mathrm{d}\theta'$$

We can get a point estimate by collapsing this posterior distribution to a single point, i.e. the best guess

# Bayes rule: example reminder

$c = \{\text{yes}, \text{no}\}$ patient having cancer?

$x \in \{-, +\}$ observed test results, a single binary feature

prior: .1% of population has cancer $p(\text{yes}) = .01$

likelihood: $p(+|\text{yes}) = .9$   TP rate of the test (90%)

$$p(c = yes \mid x) = \frac{p(c{=}yes)p(x|c{=}yes)}{p(x)}$$

FP rate of the test (5%)

posterior: $p(\text{yes}|+) = .177$

evidence: $p(+) = p(\text{yes})p(+|\text{yes}) + p(\text{no})p(+|\text{no}) = .001 \times .9 + .999 \times .05 = .05$

# Beta distribution prior

in our coin example, we know the form of likelihood:

| prior | $p(\theta)?$ |
|---|---|
| posterior | $p(\theta\|\mathcal{D})?$ |
| likelihood | $p(\mathcal{D}\|\theta) = \prod_{x \in \mathcal{D}} \mathrm{Bernoulli}(x; \theta) = \theta^{N_h}(1-\theta)^{N_t}$ |

proportional

posterior $\propto$ prior $\times$ likelihood

$$p(\theta : \alpha') \propto p(\theta : \alpha) \times p(\mathcal{D}|\theta)$$

conjugate

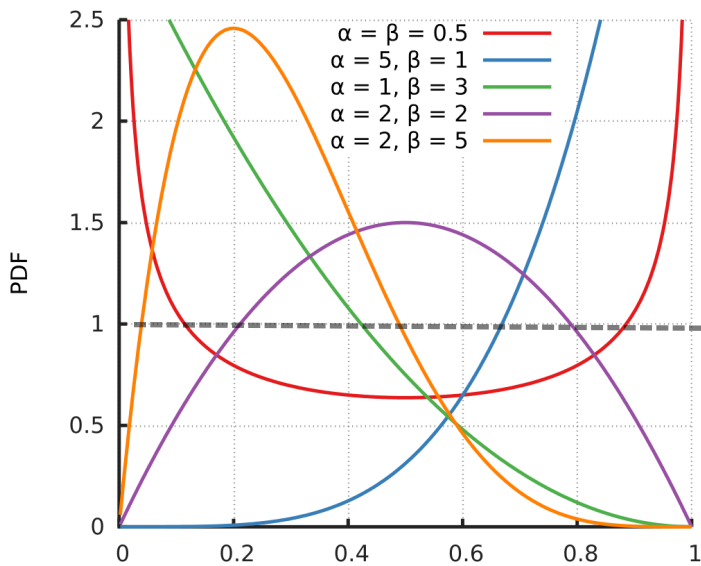A common type of prior is :  $p(\theta|a,b) \propto \theta^a (1-\theta)^b$

this means there is a normalization constant that does not depend on $\theta$

distribution of this form has a name, **Beta** distribution

we say Beta distribution is a conjugate prior to the Bernoulli likelihood

(so that we can easily update our belief with new observations, i.e. closed under Bayesian updating)

# Beta distribution

**Beta distribution** has the following density



$$\mathrm{Beta}(\theta|\alpha,\beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1}$$
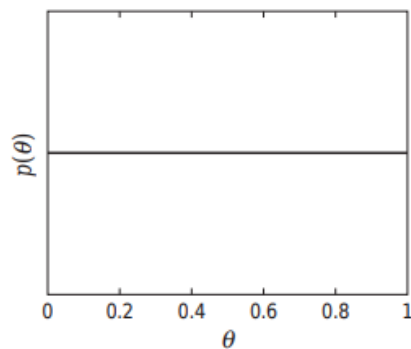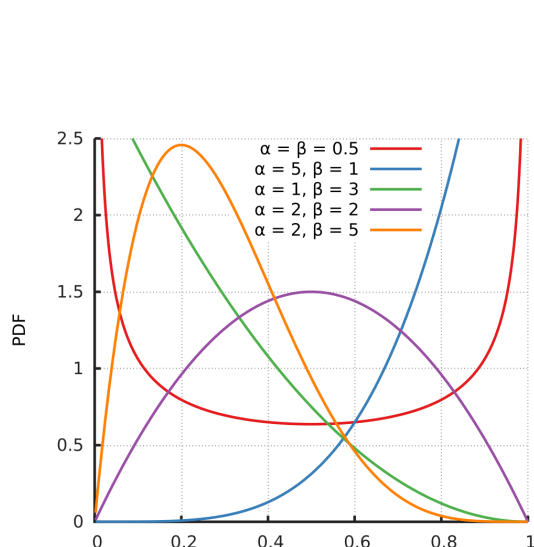
$\alpha, \beta > 0$

normalization

$\Gamma$ is the generalization of factorial to real number $\Gamma(a+1) = a\Gamma(a)$

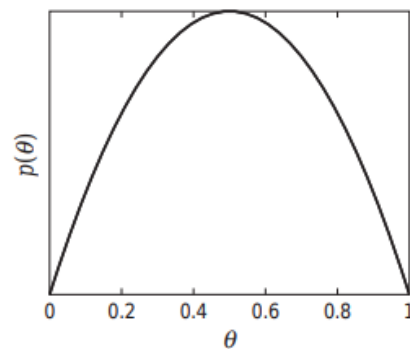$\mathrm{Beta}(\theta|\alpha=\beta=1)$ is uniform

mean of the distribution is $\mathbb{E}[\theta] = \frac{\alpha}{\alpha+\beta}$

for $\alpha, \beta > 1$ the dist. is unimodal; its mode is $\frac{\alpha-1}{\alpha+\beta-2}$
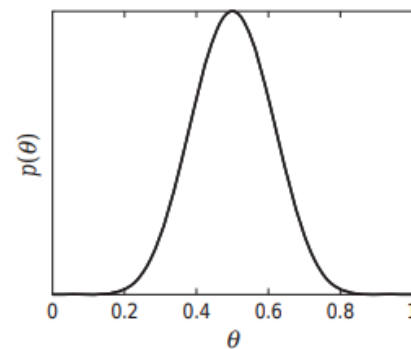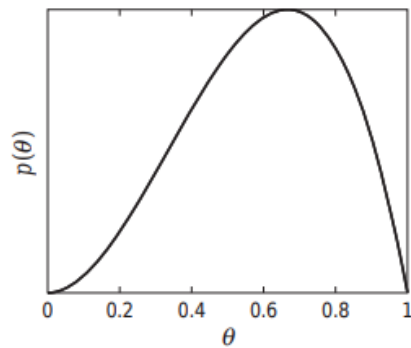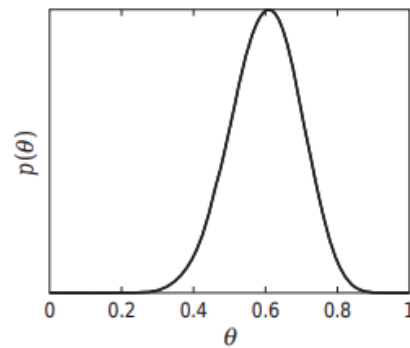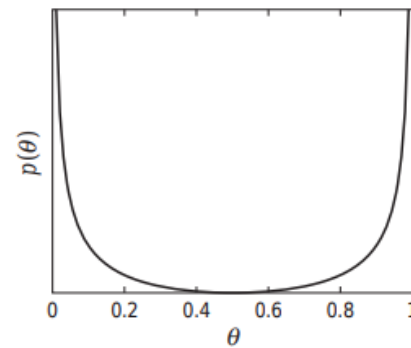
# Beta distribution: more examples

# Beta-Bernoulli posterior distribution

how to model probability of heads when we toss a coin $N$ times

proportional

| posterior | $\propto$ | prior | $\times$ | likelihood |

| prior | $p(\theta) \propto \theta^{\alpha-1}(1-\theta)^{\beta-1}$ | $p(\theta) = \text{Beta}(\theta|\alpha, \beta)$ |

| likelihood | $p(\mathcal{D}|\theta) = \theta^{N_h}(1-\theta)^{N_t}$ | $L(\theta; \mathcal{D}) = \prod \text{Bernoulli}(N_h, N_t|\theta)$ |

*product of Bernoulli likelihoods*
*equivalent to Binomial likelihood*

| posterior | $p(\theta|\mathcal{D}) \propto \theta^{\alpha+N_h-1}(1-\theta)^{\beta+N_t-1}$ | $p(\theta|\mathcal{D}) = \text{Beta}(\theta|\alpha + N_h, \beta + N_t)$ |

$\alpha, \beta$  are called *pseudo-counts*

their effect is similar to imaginary observation of heads ( $\alpha$ ) and tails ( $\beta$ )
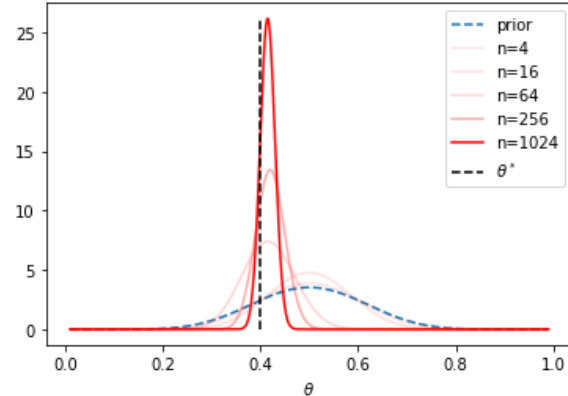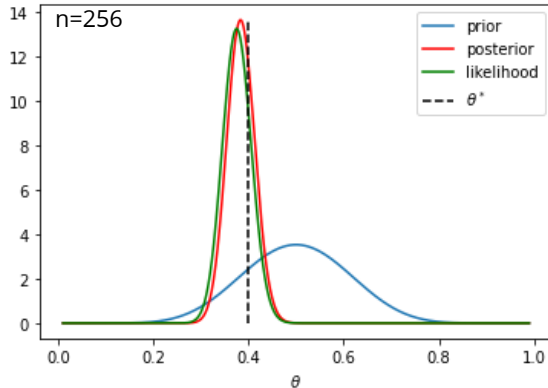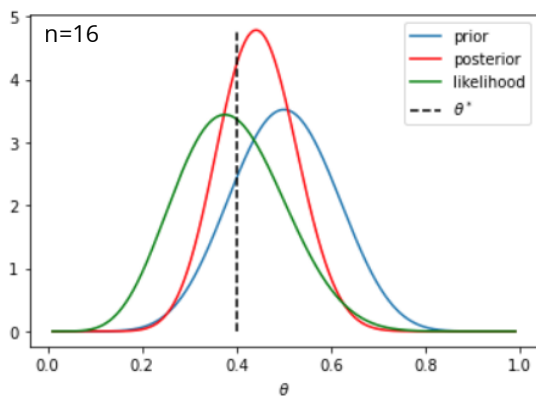
# Effect of more data

with few observations, prior has a high influence
as we increase the number of observations $N = |\mathcal{D}|$ the effect of prior diminishes
the likelihood term dominates the posterior

**example**   prior   $\text{Beta}(\theta|10, 10)$

plot of the posterior density with **n** observations

$$p(\theta|\mathcal{D}) \propto \theta^{10+H}(1-\theta)^{10+N-H}$$

# Posterior predictive

our goal was to estimate the parameters ( $\theta$ ) so that we can make predictions

what if we use the maximum likelihood estimate for the best parameter, $\theta^{MLE}$, and plug it in the $p(x|\theta)$ to make the prediction?

Example:

if we see four heads in a row, what is the probability of seeing a tail next?

if $\mathcal{D} = \{1, 1, 1, 1\}$, what is $\theta^{MLE}$? 1.0

$$\Rightarrow 1 - \theta^{MLE} = 0.0$$

$$p(0|\theta) = \theta^0 (1 - \theta)^{(1-0)} = 1 - \theta$$

Next, let's use the posterior distribution we learn through Bayesian inference

# Posterior predictive

our goal was to estimate the parameters ( $\theta$ ) so that we can make predictions

now we have a (posterior) **distribution** over parameters, $p(\theta|\mathcal{D})$, rather than a single $\theta^{MLE}$

$\theta^{MLE}$ only gives a single best guess based on that parameter, $p(x|\theta)$

To make predictions, we calculate the average prediction over all possible values of $\theta$

$$p(x|\mathcal{D}) = \int_\theta p(\theta|\mathcal{D})p(x|\theta)\mathrm{d}\theta$$

for each possible $\theta$, weight the prediction by the posterior probability of that parameter being true

posterior predictive
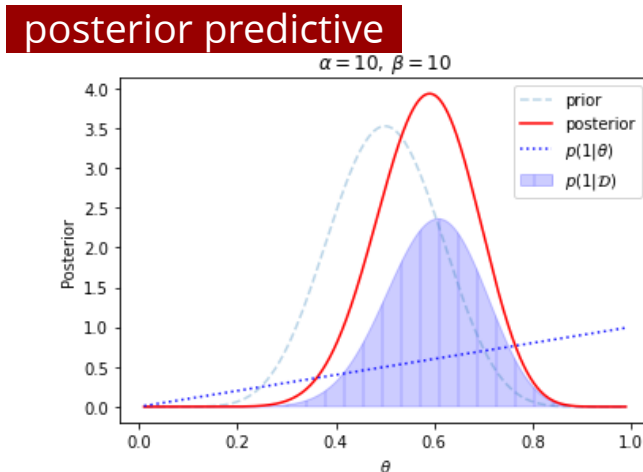


19

# Posterior predictive

our goal was to estimate the parameters ( $\theta$ ) so that we can make predictions

now we have a (posterior) **distribution** over parameters, $p(\theta|\mathcal{D})$

To make predictions, we calculate the average prediction over all possible values of $\theta$

Example   if we see four heads in a row, what is the probability of seeing a tail next?
if $\mathcal{D} = \{1, 1, 1, 1\}$, what is $p(0|\mathcal{D})$?   depends on our prior belief



when the strenght of prior gets close to zero the prediction becomes similar to MLE

# Posterior predictive for Beta-Bernoulli



start from a Beta prior  $p(\theta) = \mathrm{Beta}(\theta|\alpha, \beta)$

observe $N_h$ heads and $N_t$ tails, the posterior is  $p(\theta|\mathcal{D}) = \mathrm{Beta}(\theta|\alpha + N_h, \beta + N_t)$

Given this estimate of the parameters from training data,
how can we predict the future?

what is the probability that the next coin flip is head?

marginalize over $\theta$
$$p(x = 1|\mathcal{D}) = \int_\theta \mathrm{Bernoulli}(x = 1|\theta)\mathrm{Beta}(\theta|\alpha + N_h, \beta + N_t)\mathrm{d}\theta$$

$$= \int_\theta \theta\, \mathrm{Beta}(\theta|\alpha + N_h, \beta + N_t)d\theta = \frac{\alpha + N_h}{\alpha + \beta + N}$$

*mean of Beta dist.*



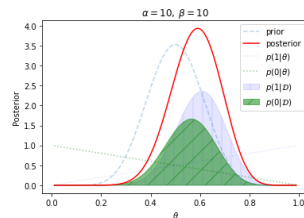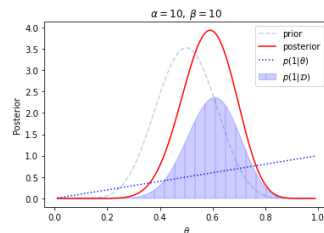**Example**   if we see four heads in a row, what is the probability of seeing a tail next?
if $\mathcal{D} = \{1, 1, 1, 1\}$, what is $p(1|\mathcal{D})$?  $\frac{14}{24}$,  $p(0|\mathcal{D})$? $\frac{10}{24}$
when we assume the prior is $\mathrm{Beta}(\alpha = 10, \beta = 10)$

compare with prediction of maximum-likelihood:  $p(x = 1|\mathcal{D}) = \frac{N_h}{N} = 1$, $p(x = 1|\mathcal{D}) = 1$

21

# Posterior predictive for Beta-Bernoulli

start from a Beta prior $p(\theta) = \text{Beta}(\theta|\alpha, \beta)$

observe $N_h$ heads and $N_t$ tails, the posterior is $p(\theta|\mathcal{D}) = \text{Beta}(\theta|\alpha + N_h, \beta + N_t)$

Given this estimate of the parameters from training data, how can we predict the future?

$$p(x = 1|\mathcal{D}) = \int_\theta \text{Bernoulli}(x = 1|\theta)\text{Beta}(\theta|\alpha + N_h, \beta + N_t)\text{d}\theta = \frac{\alpha + N_h}{\alpha + \beta + N}$$

compare with prediction of maximum-likelihood: $p(x = 1|\mathcal{D}) = \frac{N_h}{N}$

if we assume a uniform prior, the posterior predictive is $p(x = 1|\mathcal{D}) = \frac{N_h + 1}{N + 2}$

**Laplace smoothing**

a.k.a. add-one smoothing
to avoid ruling out unseen
cases with zero counts

**Example:**
sequential Baysian
updating
with uniform prior
$(N_h, N_t)$

# Strength of the prior

with a **strong prior** we need many samples to really change the posterior
for Beta distribution $\alpha + \beta$ decides how strong the prior is: how confident we are in our prior

<span style="background:#8B0000;color:white">example</span>  as our dataset grows our estimate becomes more accurate



different prior means $\frac{\alpha}{\alpha+\beta}$

posterior estimates

different prior strength $\alpha + \beta$

true value →

$p(x = 1 | \mathcal{D})$

$N$ = # samples

M = # samples

23

# Maximum a Posteriori (MAP)

sometimes it is difficult to work with the posterior dist. over parameters

**alternative**: use the parameter with the highest posterior probability $p(\theta|\mathcal{D})$

**MAP estimate** $\quad \theta^{MAP} = \arg\max_\theta p(\theta|\mathcal{D}) = \arg\max_\theta p(\theta)p(\mathcal{D}|\theta)$

compare with max-likelihood estimate  *(the only difference is in the prior term)*

$$\theta^{MLE} = \arg\max_\theta p(\mathcal{D}|\theta)$$

**example** $\quad$ for the posterior $\quad p(\theta|\mathcal{D}) = \text{Beta}(\theta|\alpha + N_h, \beta + N_t)$

MAP estimate is the **mode** of posterior $\quad \theta^{MAP} = \frac{\alpha + N_h - 1}{\alpha + \beta + N_h + N_t - 2}$

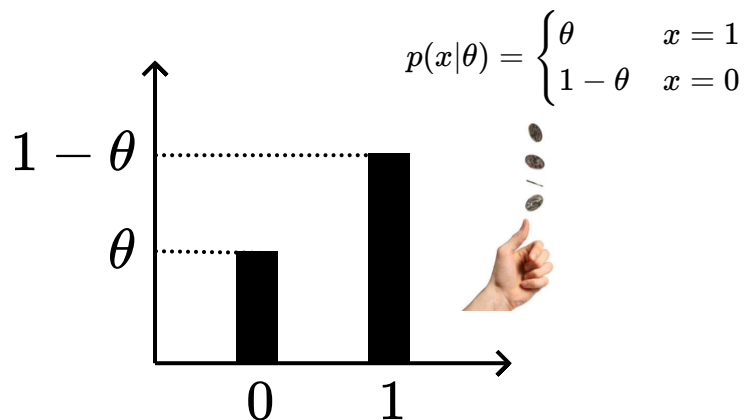compare with MLE $\quad \theta^{MLE} = \frac{N_h}{N_h + N_t}$

they are equal for uniform prior $\quad \alpha = \beta = 1$



$\mathcal{D} = \{1,1,1,1\}$
$\alpha = 10, \beta = 10$



$\mathcal{D} = \{0,0,0,1\}$
$\alpha = 10, \beta = 10$

# Categorical distribution

what if we have more than two categories (e.g., loaded dice instead of coin)

instead of Bernoulli we have multinoulli or **categorical** dist.

$$\text{Bernoulli}(x|\theta) = \theta^x (1-\theta)^{(1-x)}$$

# categories

$$\text{Cat}(x|\theta) = \prod_{k=1}^{K} \theta_k^{\mathbb{I}(x=k)}$$

$$p(x|\theta) = \begin{cases} \theta & x = 1 \\ 1-\theta & x = 0 \end{cases}$$



$1 - \theta$

$\theta$

0      1

once:              Bernoulli distribution
n times:           binomial distribution

$$p(x|\theta) = \begin{cases} \theta_1 & x = 1 \\ \theta_2 & x = 2 \\ \theta_3 & x = 3 \\ \theta_4 & x = 4 \\ \theta_5 & x = 5 \\ \theta_6 & x = 6 \end{cases}$$



$\theta_6$

$\theta_{1:5}$

1    2    3    4    5    6

categorical distribution
multinomial distribution

# Categorical distribution

what if we have more than two categories (e.g., loaded dice instead of coin)
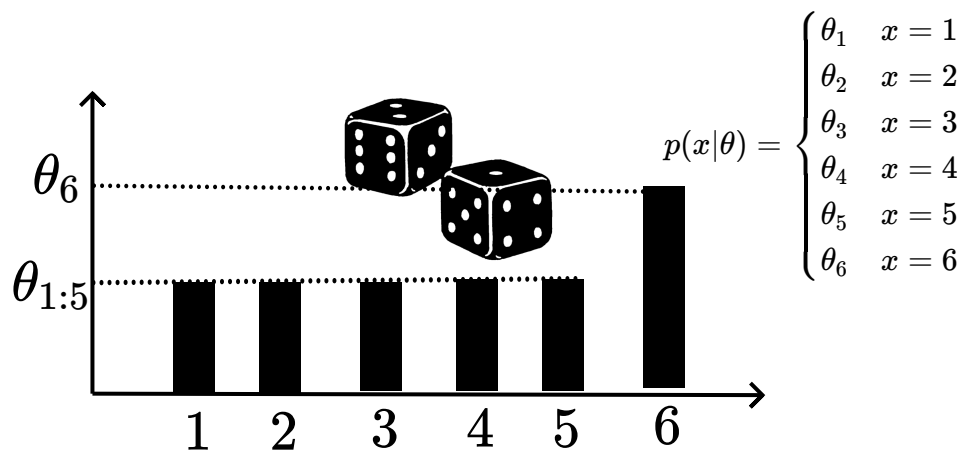instead of Bernoulli we have multinoulli or **categorical** dist.

$$\mathrm{Cat}(x|\theta) = \prod_{k=1}^{K} \theta_k^{\mathbb{I}(x=k)}$$

$$\text{where} \quad \sum_k \theta_k = 1$$

$\theta$ *belongs to probability simplex*

$$p(x|\theta) = \begin{cases} \theta_1 & x = 1 \\ \theta_2 & x = 2 \\ \theta_3 & x = 3 \\ \theta_4 & x = 4 \\ \theta_5 & x = 5 \\ \theta_6 & x = 6 \end{cases}$$

$$\sum_k^6 \theta_k = 1$$



$\theta_1 + \theta_2 + \theta_3 = 1$

$(0,0,1)$

$(0,1,0) \quad \theta_2$

$(1,0,0)$

$\theta_1$

$\theta_3$

$$K = 3$$

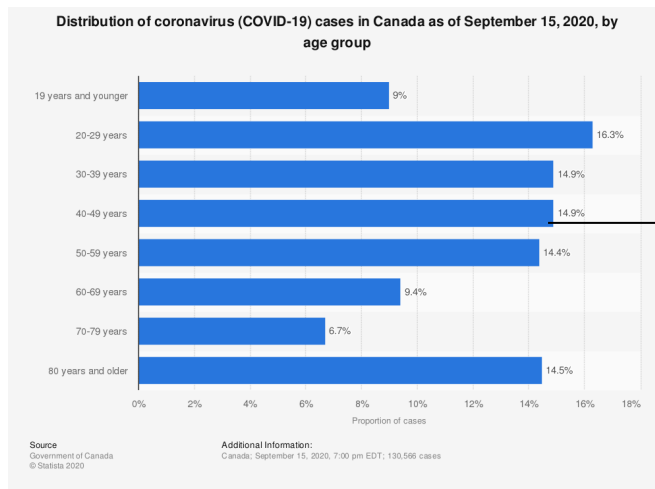# Maximum likelihood for categorical dist.

likelihood
$$p(\mathcal{D}|\theta) = \prod_{x\in\mathcal{D}} \text{Cat}(x|\theta) = \prod_{x\in\mathcal{D}} \prod_{k=1}^{K} \theta_k^{\mathbb{I}(x=k)} = \prod_{k=1}^{K} \theta_k^{N_k} \; , \; N_k = \sum_{x\in\mathcal{D}} \mathbb{I}(x=k)$$

log-likelihood
$$\ell(\theta, \mathcal{D}) = \sum_{x\in\mathcal{D}} \sum_k \mathbb{I}(x=k) \log(\theta_k) = \sum_k N_k \log(\theta_k)$$

we need to solve  $\frac{\partial}{\partial \theta_k} \ell(\theta, \mathcal{D}) = 0$  subject to  $\sum_k \theta_k = 1$    using Lagrange multipliers

similar to the binary case, max-likelihood estimate is given by data-frequencies   $\theta_k^{MLE} = \frac{N_k}{N}$
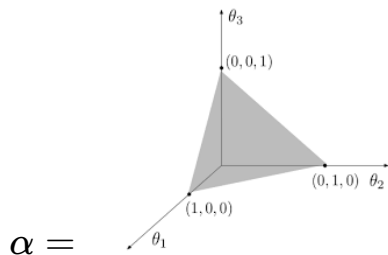


example

Distribution of coronavirus (COVID-19) cases in Canada as of September 15, 2020, by age group

| age group | proportion |
|---|---|
| 19 years and younger | 9% |
| 20-29 years | 16.3% |
| 30-39 years | 14.9% |
| 40-49 years | 14.9% |
| 50-59 years | 14.4% |
| 60-69 years | 9.4% |
| 70-79 years | 6.7% |
| 80 years and older | 14.5% |

Source
Government of Canada
© Statista 2020

Additional Information:
Canada; September 15, 2020, 7:00 pm EDT; 130,566 cases

categorical distribution with K=8

frequencies are max-likelihood parameter estimates

$\theta_5^{MLE} = .149$

# Dirichlet distribution
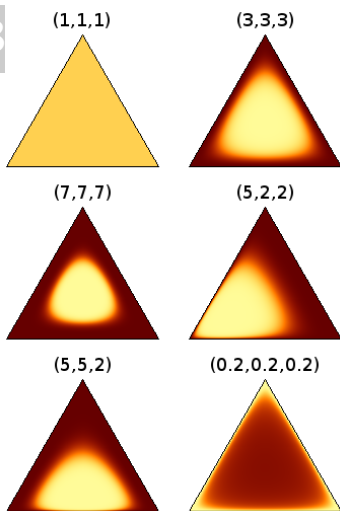


$\alpha =$

$K = 3$

is a distribution over the parameters $\theta$ of a Categorical dist.

is a generalization of Beta distribution to K categories

this should be a dist. over prob. simplex $\sum_k \theta_k = 1$

$$\mathrm{Dir}(\theta|\alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1}$$

*normalization constant*

vector of psedo-counts for K categories *(aka concentration parameters)*
$\alpha_k > 0 \; \forall k$

for $\alpha = [1, \ldots, 1]$, we get uniform distribution

for K=2, it reduces to Beta distribution

(1,1,1)   (3,3,3)

(7,7,7)   (5,2,2)

(5,5,2)   (0.2,0.2,0.2)

$\mathrm{Dir}(\theta, [.2, .2, .2])$

28

# Dirichlet-Categorical conjugate pair

Dirichlet dist. $\mathrm{Dir}(\theta|\alpha) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k - 1}$ is a conjugate prior for Categorical dist. $\mathrm{Cat}(x|\theta) = \prod_k \theta_k^{\mathbb{I}(x=k)}$

posterior $\propto$ prior $\times$ likelihood

prior $\quad p(\theta) = \mathrm{Dir}(\theta|\alpha) \propto \prod_k \theta_k^{\alpha_k - 1}$

likelihood $\quad p(\mathcal{D}|\theta) = \prod_k \theta_k^{N_k}$ $\quad$ we observe $\overset{\eta}{N_1, \ldots, N_K}$ values from each category

posterior $\quad p(\theta|\mathcal{D}) = \mathrm{Dir}(\theta|\alpha + \eta) \propto \prod_k \theta_k^{N_k + \alpha_k - 1}$ $\quad$ again, we add the real counts to pseudo-counts

posterior predictive $\quad p(x = k|\mathcal{D}) = \frac{\alpha_k + N_k}{\sum_{k'} \alpha_{k'} + N_{k'}}$

MAP $\quad \theta_k^{MAP} = \frac{\alpha_k + N_k - 1}{(\sum_{k'} \alpha_{k'} + N_{k'}) - K}$

# Summary

in ML we often build a probabilistic model of the data $p(x; \theta)$

learning a good model could mean **maximizing the likelihood** of the data

$$\max_\theta \log p(\mathcal{D}|\theta) \Big|$$ sometimes closed form solution
for more complex p, we use numerical methods

an alternative is a **Bayesian approach**:

- maintain a **distribution** over model parameters
- can specify our **prior** knowledge $p(\theta)$
- we can use **Bayes rule** to update our belief after new oabservation $p(\theta|\mathcal{D})$
- we can make predictions using **posterior predictive** $p(x|\mathcal{D})$
- can be computationally **expensive** *(not in our examples so far)*

a middle path is **MAP estimate**: $\max_\theta \log p(\mathcal{D}|\theta)p(\theta)$

- models our **prior** belief
- use a single point estimate and picks the model with highest posterior probability