

Introduction to Bayesian methods

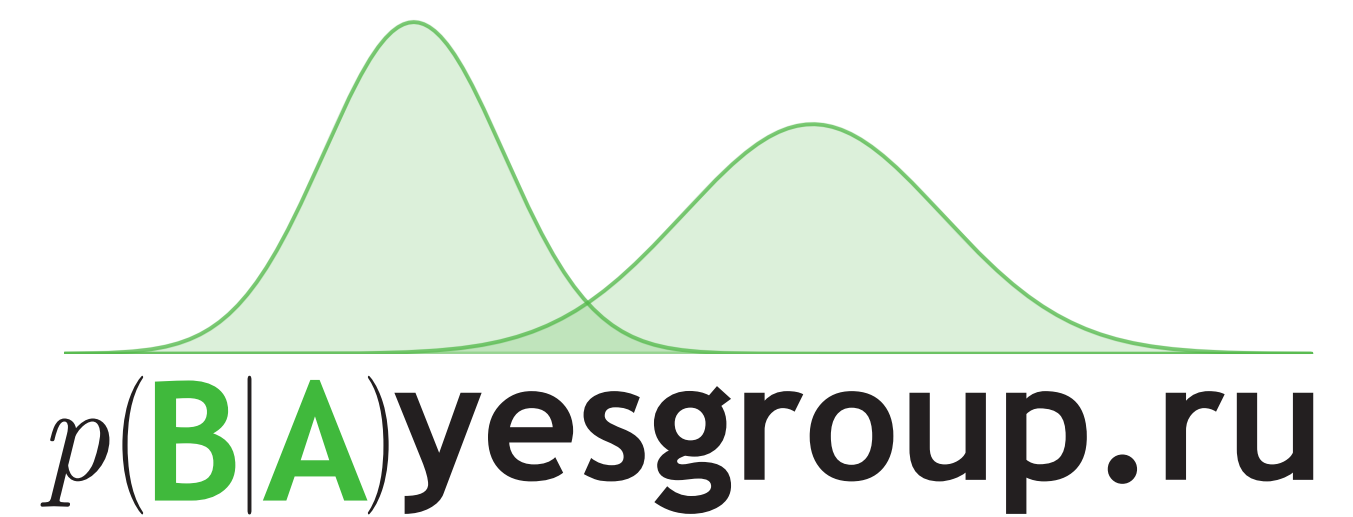
Ekaterina Lobacheva

Higher School of Economics, Samsung-HSE Laboratory
Moscow, Russia



NATIONAL RESEARCH
UNIVERSITY

SAMSUNG
Research



Problem set

The problem set is
available here:

tiny.cc/ASGM_bayes_problems

Problem 1: Bayesian reasoning

Setting

During medical checkup, one of the tests indicates a serious disease. The test has high accuracy 99% (probability of true positive is 99%, probability of true negative is 99%). However, the disease is quite rare, and only one person in 10000 is affected.

Question

Calculate the probability that the examined person has the disease.

Problem 1: Bayesian reasoning

- $d \in \{0, 1\}$ — disease (1 means that the person has a disease)
- $t \in \{0, 1\}$ — test (1 means that test says that the person has a disease)

Setting: $p(t = 1 \mid d = 1) = p(t = 0 \mid d = 0) = 0.99$, $p(d = 1) = 10^{-4}$

Question: $p(d = 1 \mid t = 1) = ?$

Problem 1: Bayesian reasoning

- $d \in \{0, 1\}$ — disease (1 means that the person has a disease)
- $t \in \{0, 1\}$ — test (1 means that test says that the person has a disease)

Setting: $p(t = 1 \mid d = 1) = p(t = 0 \mid d = 0) = 0.99$, $p(d = 1) = 10^{-4}$

Question: $p(d = 1 \mid t = 1) = ?$

$$\begin{aligned} p(d = 1 \mid t = 1) &= \frac{p(t = 1 \mid d = 1)p(d = 1)}{p(t = 1 \mid d = 1)p(d = 1) + p(t = 1 \mid d = 0)p(d = 0)} = \\ &= \frac{0.99 \cdot 10^{-4}}{0.99 \cdot 10^{-4} + 0.01 \cdot (1 - 10^{-4})} \approx 1\% \end{aligned}$$

Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability θ of landing heads up
- Data: $X = (x_1, \dots, x_n)$, $x \in \{0, 1\}$

Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$



Head (H)



Tail (T)

Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability θ of landing heads up
- Data: $X = (x_1, \dots, x_n)$, $x \in \{0, 1\}$

Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$

Likelihood: $Bern(x | \theta) = \theta^x (1 - \theta)^{1-x}$



Head (H)



Tail (T)

Example: coin tossing

- We have a coin which may be fair or not
- The task is to estimate a probability θ of landing heads up
- Data: $X = (x_1, \dots, x_n)$, $x \in \{0, 1\}$



Head (H)



Tail (T)

Probabilistic model:

$$p(x, \theta) = p(x | \theta)p(\theta)$$

Likelihood: $Bern(x | \theta) = \theta^x (1 - \theta)^{1-x}$

Prior: ???

Example: coin tossing

How to choose a prior?

- Correct domain: $\theta \in [0, 1]$
- Include prior knowledge: a coin is most likely fair
- Inference complexity: use conjugate prior

Example: coin tossing

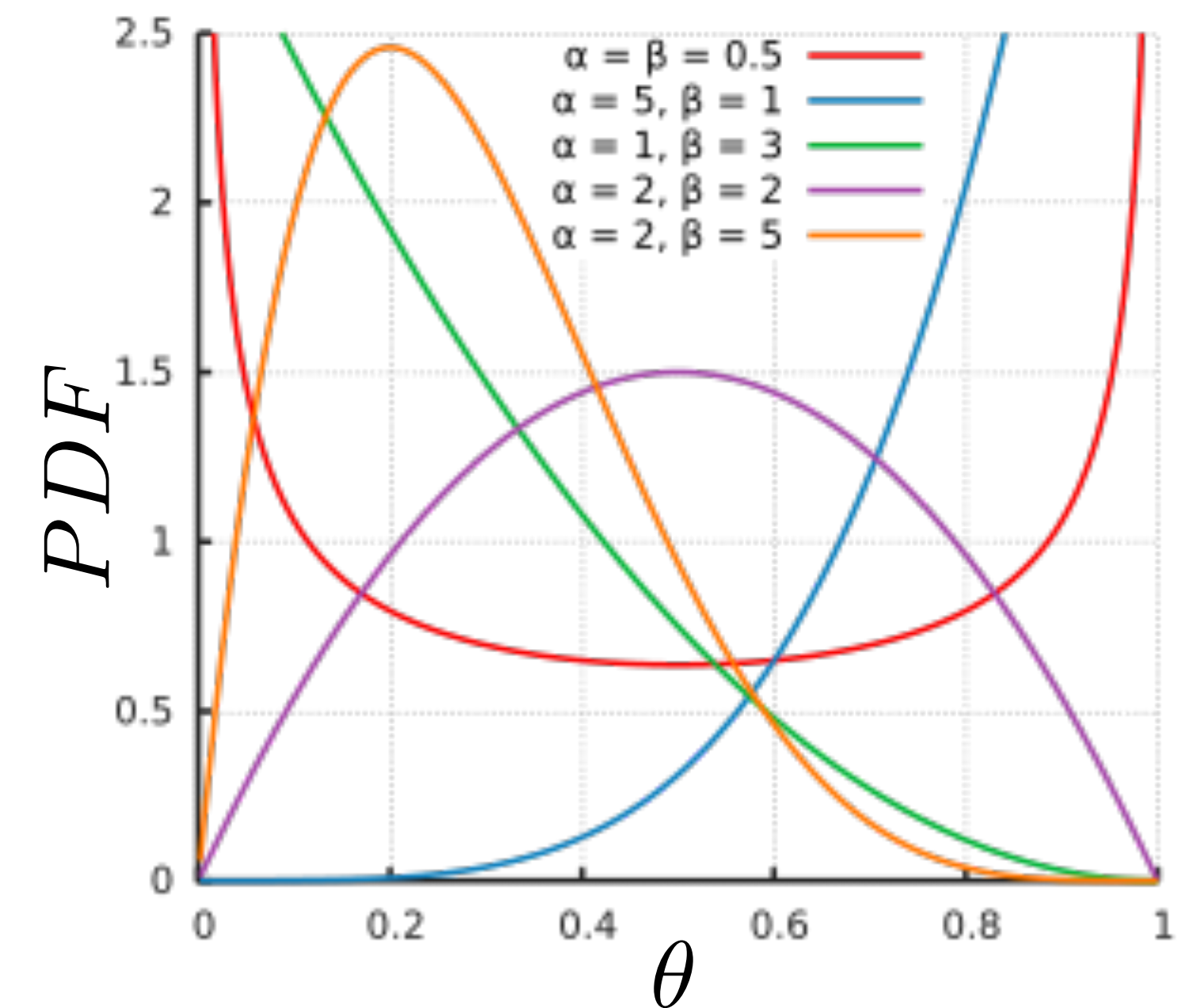
How to choose a prior?

- Correct domain: $\theta \in [0, 1]$
- Include prior knowledge: a coin is most likely fair
- Inference complexity: use conjugate prior

Beta distribution matches all requirements:

$$Beta(\theta \mid a, b) = \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Beta distribution



Example: coin tossing

Let's check that our likelihood and prior are conjugate:

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \quad p(\theta) = \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Idea — check that prior and posterior lay in the same parametric family:

Here different constants are denoted with the same letter C for demonstration reasons.

Example: coin tossing

Let's check that our likelihood and prior are conjugate:

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \quad p(\theta) = \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Idea — check that prior and posterior lay in the same parametric family:

$$p(\theta) = C \theta^C (1 - \theta)^C$$

Here different constants are denoted with the same letter C for demonstration reasons.

Example: coin tossing

Let's check that our likelihood and prior are conjugate:

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \quad p(\theta) = \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Idea — check that prior and posterior lay in the same parametric family:

$$p(\theta) = C \theta^C (1 - \theta)^C$$

$$\begin{aligned} p(\theta | x) &= \frac{1}{C} p(x | \theta) p(\theta) = \frac{1}{C} \theta^x (1 - \theta)^{1-x} \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1} = \\ &= C \theta^C (1 - \theta)^C \end{aligned}$$

Here different constants are denoted with the same letter C for demonstration reasons.

Example: coin tossing

Let's check that our likelihood and prior are conjugate:

$$p(x | \theta) = \theta^x (1 - \theta)^{1-x} \quad p(\theta) = \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1}$$

Idea — check that prior and posterior lay in the same parametric family:

$$p(\theta) = \boxed{C \theta^C (1 - \theta)^C} \text{ conjugacy}$$

$$p(\theta | x) = \frac{1}{C} p(x | \theta) p(\theta) = \frac{1}{C} \theta^x (1 - \theta)^{1-x} \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1} = \\ = \boxed{C \theta^C (1 - \theta)^C} \text{ conjugacy}$$

Here different constants are denoted with the same letter C for demonstration reasons.

Example: coin tossing

Bayesian inference after receiving data $X = (x_1, \dots, x_n)$:

$$p(\theta | X) = \frac{1}{Z} p(X | \theta) p(\theta) = \frac{1}{Z} \left[\prod_{i=1}^n p(x_i | \theta) \right] p(\theta) =$$

Example: coin tossing

Bayesian inference after receiving data $X = (x_1, \dots, x_n)$:

$$\begin{aligned} p(\theta | X) &= \frac{1}{Z} p(X | \theta) p(\theta) = \frac{1}{Z} \left[\prod_{i=1}^n p(x_i | \theta) \right] p(\theta) = \\ &= \frac{1}{Z} \left[\prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \right] \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1} = \end{aligned}$$

Example: coin tossing

Bayesian inference after receiving data $X = (x_1, \dots, x_n)$:

$$\begin{aligned} p(\theta | X) &= \frac{1}{Z} p(X | \theta) p(\theta) = \frac{1}{Z} \left[\prod_{i=1}^n p(x_i | \theta) \right] p(\theta) = \\ &= \frac{1}{Z} \left[\prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \right] \frac{1}{\text{B}(a, b)} \theta^{a-1} (1 - \theta)^{b-1} = \\ &= \frac{1}{Z'} \theta^{a + \sum_{i=1}^n x_i - 1} (1 - \theta)^{b + n - \sum_{i=1}^n x_i - 1} \end{aligned}$$

Example: coin tossing

Bayesian inference after receiving data $X = (x_1, \dots, x_n)$:

$$\begin{aligned} p(\theta | X) &= \frac{1}{Z} p(X | \theta) p(\theta) = \frac{1}{Z} \left[\prod_{i=1}^n p(x_i | \theta) \right] p(\theta) = \\ &= \frac{1}{Z} \left[\prod_{i=1}^n \theta^{x_i} (1 - \theta)^{1-x_i} \right] \frac{1}{B(a, b)} \theta^{a-1} (1 - \theta)^{b-1} = \\ &= \frac{1}{Z'} \theta^{a + \sum_{i=1}^n x_i - 1} (1 - \theta)^{b + n - \sum_{i=1}^n x_i - 1} = \text{Beta}(\theta | a', b') \end{aligned}$$

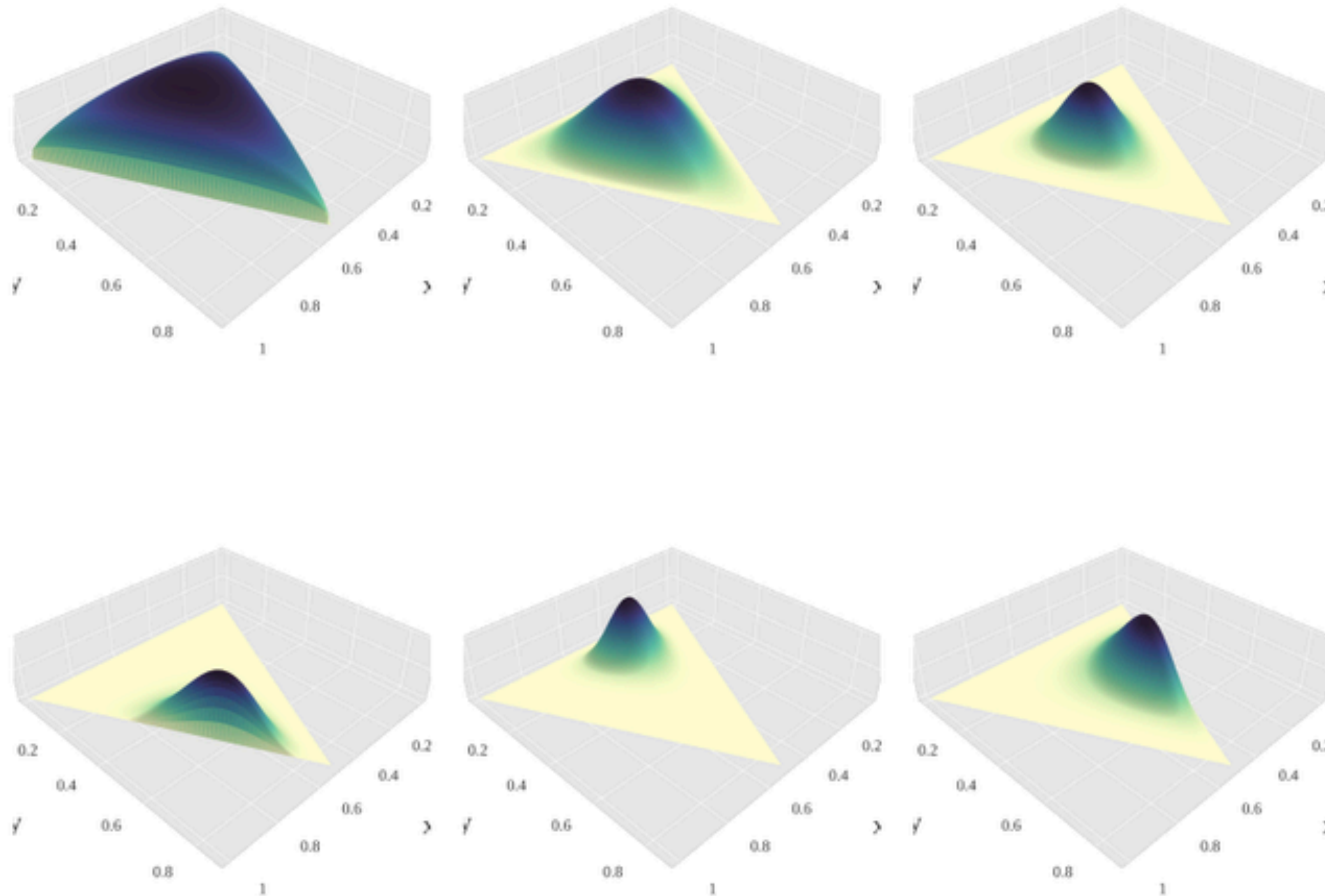
New parameters: $a' = a + \sum_{i=1}^n x_i$ $b' = b + n - \sum_{i=1}^n x_i$

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in \mathcal{S}_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Dirichlet distribution



Beta distribution is a special case of Dirichlet distribution:

$$\text{Dir}(\theta \mid \alpha) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$\text{Beta}(\theta \mid a, b) \propto \theta^{a-1} (1 - \theta)^{b-1}$$

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Questions

- Check that likelihood and prior are conjugate
- Compute the posterior $p(\theta | X, \alpha)$
- * Compare $\mathbb{E}_{p(\theta | X, \alpha)} \theta$ and θ_{ML}
- * Compute the predictive posterior $p(x_{N+1} = j | X, \alpha)$

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Questions

- Check that likelihood and prior are conjugate
- Compute the posterior $p(\theta | X, \alpha)$
- * Compare $\mathbb{E}_{p(\theta | X, \alpha)} \theta$ and θ_{ML}
- * Compute the predictive posterior $p(x_{N+1} = j | X, \alpha)$

Problem 2: Bayesian framework

Probabilistic model: $p(X, \theta) = p(X | \theta)p(\theta) = p(X | \theta)Dir(\theta | \alpha)$

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$Dir(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Here different constants are denoted with the same letter C for demonstration reasons.

Problem 2: Bayesian framework

Probabilistic model: $p(X, \theta) = p(X | \theta)p(\theta) = p(X | \theta)Dir(\theta | \alpha)$

Prior:
$$p(\theta) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} = C \prod_{k=1}^K \theta_k^C$$

Here different constants are denoted with the same letter C for demonstration reasons.

Problem 2: Bayesian framework

Probabilistic model: $p(X, \theta) = p(X | \theta)p(\theta) = p(X | \theta)Dir(\theta | \alpha)$

Prior:
$$p(\theta) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} = C \prod_{k=1}^K \theta_k^C$$

Posterior:
$$p(\theta | X) \propto p(X | \theta)p(\theta) = \prod_{k=1}^K \theta_k^{N_k} \cdot \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} =$$
$$= C \prod_{k=1}^K \theta_k^C$$

Here different constants are denoted with the same letter C for demonstration reasons.

Problem 2: Bayesian framework

Probabilistic model: $p(X, \theta) = p(X | \theta)p(\theta) = p(X | \theta)Dir(\theta | \alpha)$

Prior:
$$p(\theta) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} = C \prod_{k=1}^K \theta_k^C$$

Posterior:
$$p(\theta | X) \propto p(X | \theta)p(\theta) = \prod_{k=1}^K \theta_k^{N_k} \cdot \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} =$$

$$= C \prod_{k=1}^K \theta_k^C$$

conjugate

Here different constants are denoted with the same letter C for demonstration reasons.

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Questions

- Check that likelihood and prior are conjugate
- Compute the posterior $p(\theta | X, \alpha)$
- * Compare $\mathbb{E}_{p(\theta | X, \alpha)} \theta$ and θ_{ML}
- * Compute the predictive posterior $p(x_{N+1} = j | X, \alpha)$

Problem 2: Bayesian framework

Likelihood and prior are conjugate \rightarrow posterior is Dirichlet

$$\begin{aligned} p(\theta | X) &\propto p(X | \theta)p(\theta) = \prod_{k=1}^K \theta_k^{N_k} \cdot \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \propto \\ &\propto \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} \end{aligned}$$

Problem 2: Bayesian framework

Likelihood and prior are conjugate \rightarrow posterior is Dirichlet

$$\begin{aligned} p(\theta | X) &\propto p(X | \theta)p(\theta) = \prod_{k=1}^K \theta_k^{N_k} \cdot \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1} \propto \\ &\propto \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} \end{aligned}$$

$$p(\theta | X) = \text{Dir}(\theta | \alpha'), \quad \alpha' = (\alpha_1 + N_1, \dots, \alpha_K + N_K)$$

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Questions

- Check that likelihood and prior are conjugate
- Compute the posterior $p(\theta | X, \alpha)$
- * Compare $\mathbb{E}_{p(\theta | X, \alpha)} \theta$ and θ_{ML}
- * Compute the predictive posterior $p(x_{N+1} = j | X, \alpha)$

Problem 2: frequentist framework

θ is restricted to simplex. To omit the inequality restrictions change parameterization to $\mu_k = \log \theta_k$, $\mu_k \in \mathbb{R}$

The Lagrangian has the form:

$$\begin{aligned}\mathcal{L}(\mu, \lambda) &= \log p(X \mid \exp \mu) - \lambda(\sum_{k=1}^K \exp \mu_k - 1) = \\ &= \sum_{k=1}^K (N_k \mu_k - \lambda \exp \mu_k) + \lambda\end{aligned}$$

Differentiation:

$$\begin{aligned}0 &= \frac{\partial \mathcal{L}(\mu, \lambda)}{\partial \mu_k} = N_k - \lambda \exp \mu_k \Rightarrow \theta_k = \exp \mu_k = \frac{N_k}{\lambda} \\ 0 &= \frac{\partial \mathcal{L}(\mu, \lambda)}{\partial \lambda} = -\sum_{k=1}^K \exp \mu_k + 1 \Rightarrow \lambda = \sum_{k=1}^K N_k\end{aligned} \quad \rightarrow \quad \theta_k = \frac{N_k}{\sum_{l=1}^K N_l}$$

Problem 2: Bayesian framework

Maximum likelihood estimate: $\theta_k = \frac{N_k}{\sum_{l=1}^K N_l}$

Expectation of the posterior: $\mathbb{E}_{p(\theta|X)} \theta_k = \frac{\alpha_k + N_k}{\sum_{l=1}^K \alpha_l + N_l}$

Small K \rightarrow Bayesian estimate is mostly based on prior

Large K \rightarrow Bayesian estimate is very similar to ML estimate

Problem 2: Bayesian framework

Setting

- $p(X | \theta) = \prod_{k=1}^K \theta_k^{N_k}$ — multinomial likelihood, $\theta \in S_K$
- Dirichlet prior:
$$\text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha_1, \dots, \alpha_K)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Questions

- Check that likelihood and prior are conjugate
- Compute the posterior $p(\theta | X, \alpha)$
- * Compare $\mathbb{E}_{p(\theta | X, \alpha)} \theta$ and θ_{ML}
- * Compute the predictive posterior $p(x_{N+1} = j | X, \alpha)$

Problem 2: Bayesian framework

$$p(x_{N+1} = j \mid X, \alpha) = \int_{S_K} p(x_{N+1} = j \mid \theta) p(\theta \mid X, \alpha) d\theta =$$

Problem 2: Bayesian framework

$$\begin{aligned} p(x_{N+1} = j \mid X, \alpha) &= \int_{S_K} p(x_{N+1} = j \mid \theta) p(\theta \mid X, \alpha) d\theta = \\ &= \frac{\int_{S_K} \theta_j \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1} d\theta}{B(\alpha_1 + N_1, \dots, \alpha_K + N_K)} = \frac{B(\alpha_1 + N_1, \dots, \alpha_j + N_j + 1, \dots, \alpha_K + N_K)}{B(\alpha_1 + N_1, \dots, \alpha_j + N_j, \dots, \alpha_K + N_K)} = \\ &= \frac{\Gamma(\alpha_1 + N_1) \dots \Gamma(\alpha_j + N_j + 1) \dots \Gamma(\alpha_K + N_K)}{\Gamma(\alpha_1 + N_1) \dots \Gamma(\alpha_j + N_j) \dots \Gamma(\alpha_K + N_K)} \cdot \frac{\Gamma(\sum_l (\alpha_l + N_l))}{\Gamma(\sum_l (\alpha_l + N_l) + 1)} = \\ &= \frac{\alpha_j + N_j}{\sum_k \alpha_k + N} \end{aligned}$$