

CS 2750 Machine Learning  
Lecture 6

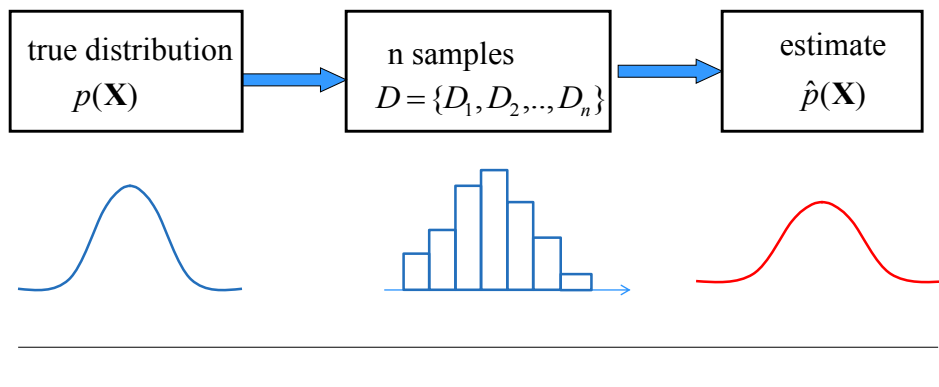
Density estimation II

Milos Hauskrecht  
[milos@pitt.edu](mailto:milos@pitt.edu)  
5329 Sennott Square

Density estimation

**Data:**  $D = \{D_1, D_2, \dots, D_n\}$   
 $D_i = \mathbf{x}_i$  a vector of attribute values

**Objective:** estimate the model of the underlying probability distribution over variables  $\mathbf{X}$ ,  $p(\mathbf{X})$ , using examples in  $D$



## Density estimation

Types of density estimation:

### Parametric

- the distribution is modeled using a set of parameters  $\Theta$   
$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta)$$
- **Example:** mean and covariances of a multivariate normal
- **Estimation:** find parameters  $\Theta$  describing data  $D$

### Non-parametric

- The model of the distribution utilizes all examples in  $D$
- As if all examples were parameters of the distribution
- **Examples:** Nearest-neighbor

## ML Parameter estimation

**Model**  $\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta_{ML})$     **Data**  $D = \{D_1, D_2, \dots, D_n\}$

- **Maximum likelihood (ML)**

$$\Theta_{ML} = \arg \max_{\Theta} P(D | \Theta, \xi)$$

$$\begin{aligned} p(D | \Theta, \xi) &= P(D_1, D_2, \dots, D_n | \Theta, \xi) \\ &= P(D_1 | \Theta, \xi) P(D_2 | \Theta, \xi) \dots P(D_n | \Theta, \xi) \\ &= \prod_{i=1}^n P(D_i | \Theta, \xi) \end{aligned}$$

↓ Independent examples

**Log likelihood – has the same maximum as likelihood**

$$\log p(D | \Theta, \xi) = \sum_{i=1}^n \log P(D_i | \Theta, \xi)$$

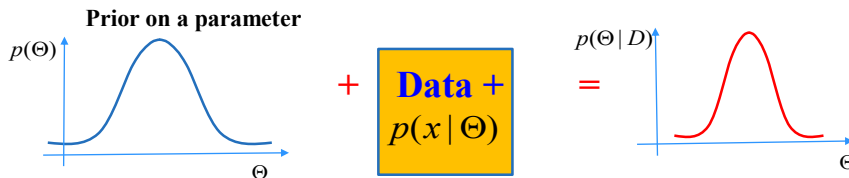
## Bayesian parameter estimation

### Bayesian parameter estimation

- Uses the posterior distribution for parameters
- Posterior ‘covers’ all possible parameter values (& their “weights”)

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{p(D | \xi)}$$

← Data Likelihood      ← Parameter prior



**Conjugate choices:** Prior distribution **matches** the data distribution → Posterior is the same type

## Parameter estimation

### Other criteria:

- **Maximum a posteriori probability (MAP)**  
(mode of the posterior)

$$\text{Model: } \hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta_{MAP})$$

$$\Theta_{MAP} = \arg \max_{\Theta} p(\Theta | D, \xi)$$

- **Expected value of the parameter**

$$\hat{p}(\mathbf{X}) = p(\mathbf{X} | \Theta_{EV})$$

- Expectation taken with regard to posterior
- (mean of the posterior)

$$\Theta_{EV} = E_{p(\Theta | D, \xi)}(\Theta) = \int \Theta p(\Theta | D, \xi) d\Theta$$

## Bernoulli distribution

### Model for random variable with two outcomes

- **Random variable:**  $x$
- **Two outcomes:** 0 or 1
- **Distribution:**

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

where  $\theta$  is the probability of  $x=1$

### Example: Coin toss

#### Outcomes:

- Head  $\rightarrow x=1$
- Tail  $\rightarrow x=0$
- $\theta \rightarrow$  probability of a Head



## Prior distribution

### Choice of prior: Beta distribution

$$p(\theta | \xi) = \text{Beta}(\theta | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1-\theta)^{\alpha_2-1}$$

$\Gamma(x)$  - a Gamma function  $\Gamma(x) = (x-1)\Gamma(x-1)$

For integer values of  $x$   $\Gamma(n) = (n-1)!$

### Why to use Beta distribution?

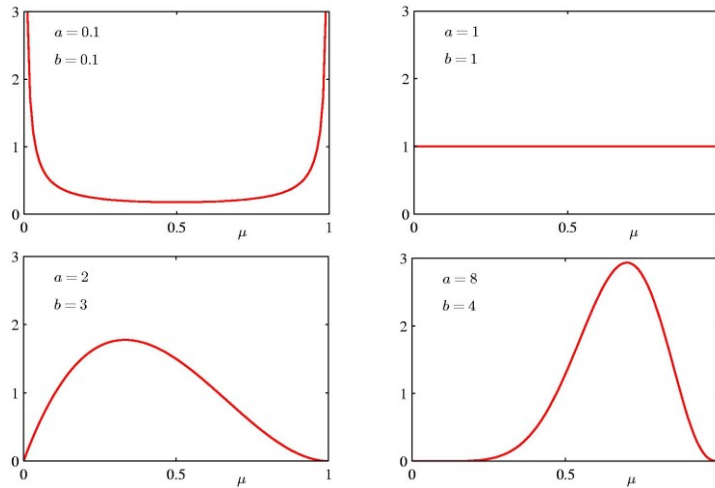
Beta distribution “fits” Bernoulli sample - **conjugate choices**

$$P(D | \theta, \xi) = \theta^{N_1} (1-\theta)^{N_2}$$

### Posterior distribution is again a Beta distribution

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

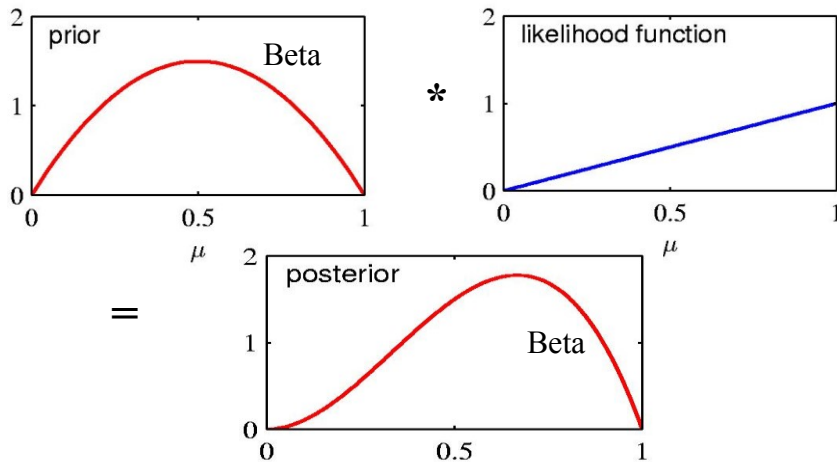
## Beta distribution



$$p(\theta | \xi) = \text{Beta}(\theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

CS 2750 Machine Learning

## Posterior distribution



$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) \text{Beta}(\theta | \alpha_1, \alpha_2)}{P(D | \xi)} = \text{Beta}(\theta | \alpha_1 + N_1, \alpha_2 + N_2)$$

## Binomial distribution



**Example problem:**  $N$  coin flips, where each coin flip can have two results: head or tail

**Outcome:**  $N_1$  - number of heads seen     $N_2$  - number of tails seen  
in  $N$  trials

**Model:** probability of a head  $\theta$   
probability of a tail  $(1-\theta)$

**Probability of an outcome:**

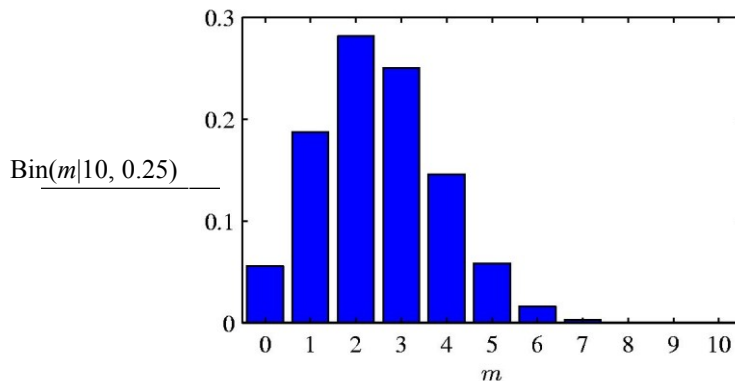
$$P(N_1 | N, \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N-N_1} \quad \text{Binomial distribution}$$

**Binomial distribution:**

- models order independent sequence of Bernoulli trials

## Binomial distribution

**Binomial distribution:**



**Matching prior: beta distribution**

## Multinomial distribution



**Example:** multiple rolls of a die with 6 results

**Outcome:** counts of occurrences of  $k$  possible outcomes of  $N$  trials:  $N_i$  - a number of times an outcome  $i$  has been seen

$$\sum_{i=1}^k N_i = N$$

**Model parameters:**  $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_k)$  s.t.  $\sum_{i=1}^k \theta_i = 1$   
 $\theta_i$  - probability of an outcome  $i$

**Probability distribution:**

$$P(N_1, N_2, \dots, N_k | \boldsymbol{\theta}, \xi) = \frac{N!}{N_1! N_2! \dots N_k!} \theta_1^{N_1} \theta_2^{N_2} \dots \theta_k^{N_k} \quad \text{Multinomial distribution}$$

**ML estimate:**

$$\theta_{i,ML} = \frac{N_i}{N}$$

## Posterior and MAP estimate



**Choice of the prior: Dirichlet distribution**

$$Dir(\boldsymbol{\theta} | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_k^{\alpha_k-1}$$

**Dirichlet is the conjugate choice for the multinomial sampling**

$$P(D | \boldsymbol{\theta}, \xi) = P(N_1, N_2, \dots, N_k | \boldsymbol{\theta}, \xi) = \frac{N!}{N_1! N_2! \dots N_k!} \theta_1^{N_1} \theta_2^{N_2} \dots \theta_k^{N_k}$$

**Posterior density**

$$p(\boldsymbol{\theta} | D, \xi) = \frac{P(D | \boldsymbol{\theta}, \xi) Dir(\boldsymbol{\theta} | \alpha_1, \alpha_2, \dots, \alpha_k)}{P(D | \xi)} = Dir(\boldsymbol{\theta} | \alpha_1 + N_1, \dots, \alpha_k + N_k)$$

**MAP estimate:**

$$\theta_{i,MAP} = \frac{\alpha_i + N_i - 1}{\sum_{i=1, \dots, k} (\alpha_i + N_i) - k}$$

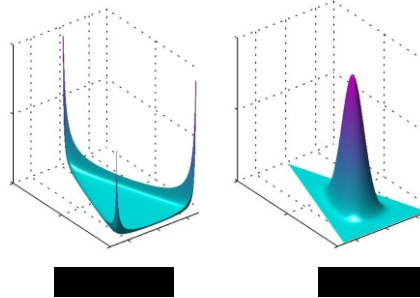
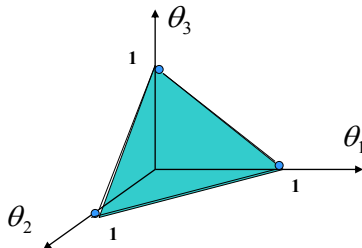
## Dirichlet distribution



**Dirichlet distribution:**

$$Dir(\boldsymbol{\theta} | \alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \theta_2^{\alpha_2-1} \dots \theta_k^{\alpha_k-1}$$

**Assume: k=3**



## Other distributions

**The same ideas can be applied to other distributions**

- Typically we choose distributions that behave well so that computations lead to “nice” solutions

- **Exponential family of distributions**

**Conjugate choices** for some of the distributions from the exponential family:

- **Binomial – Beta**
- **Multinomial - Dirichlet**
- **Exponential – Gamma**
- **Poisson – Inverse Gamma**
- **Gaussian - Gaussian (mean) and Wishart (covariance)**

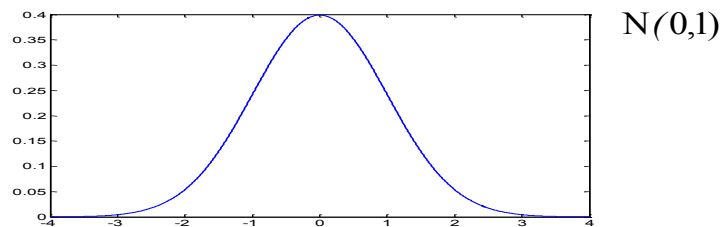


## Gaussian (normal) distribution

- **Gaussian:**  $x \sim N(\mu, \sigma)$
- **Parameters:**  $\mu$  - mean  
 $\sigma$  - standard deviation
- **Density function:**

$$p(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]$$

- **Example:**



## Parameter estimates

- **Loglikelihood**  $l(D, \mu, \sigma) = \log \prod_{i=1}^n p(x_i | \mu, \sigma)$

- **ML estimates of the mean and variance:**

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

– ML variance estimate is biased

$$E_n(\hat{\sigma}^2) = E_n\left(\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2\right) = \frac{n-1}{n} \sigma^2 \neq \sigma^2$$

- **Unbiased estimate:**

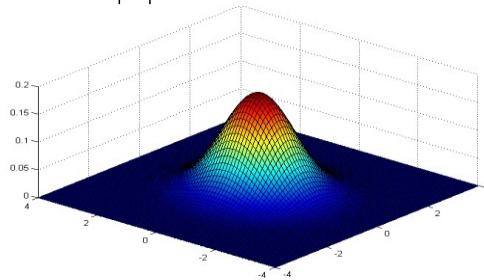
$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

## Multivariate normal distribution

- **Multivariate normal:**  $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- **Parameters:**  $\boldsymbol{\mu}$ - mean  
 $\boldsymbol{\Sigma}$ - covariance matrix
- **Density function:**

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- **Example:**

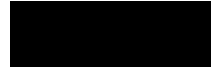


## Partitioned Gaussian Distributions

- **Multivariate Gaussian:**



- **Example:**



Precision matrix

- **What are the distributions for marginals and conditionals?**

$$p(x_a)$$

$$p(x_a | x_b)$$

## Conditionals and Marginals

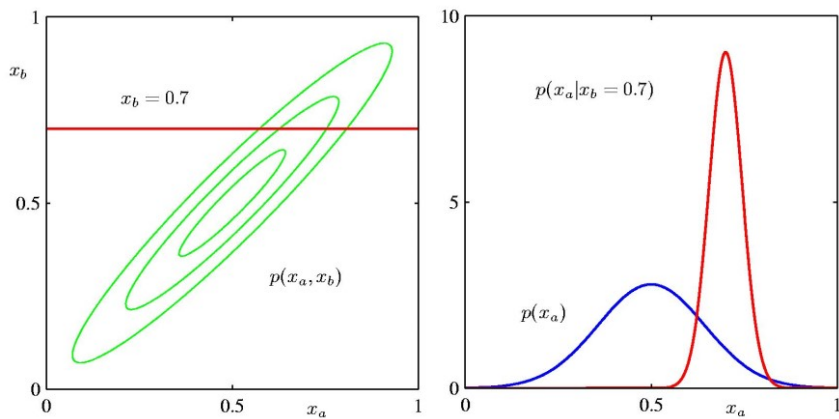
- Conditional density:



- Marginal Density:



## Conditionals and Marginals



## Parameter estimates

- **Loglikelihood**  $l(D, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

- **ML estimates of the mean and covariances:**

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

– Covariance estimate is biased

$$E_n(\hat{\boldsymbol{\Sigma}}) = E_n\left(\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T\right) = \frac{n-1}{n} \boldsymbol{\Sigma} \neq \boldsymbol{\Sigma}$$

- **Unbiased estimate:**

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

## Posterior of the mean of a multivariate normal

- **Assume a prior on the mean  $\boldsymbol{\mu}$  that is normally distributed:**

$$p(\boldsymbol{\mu}) \approx N(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$$

- **Then the posterior of  $\boldsymbol{\mu}$  is normally distributed**

$$\begin{aligned} p(\boldsymbol{\mu} | D) &\approx \left( \prod_{i=1}^n \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_i - \boldsymbol{\mu})\right] \right) \\ &\quad * \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_p|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_p)^T \boldsymbol{\Sigma}_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p)\right] \\ &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_n|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)\right] \end{aligned}$$

## Posterior of the mean of a multivariate normal

- Then the posterior of  $\boldsymbol{\mu}$  is normally distributed

$$p(\boldsymbol{\mu} | D) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_n|^{1/2}} \exp\left[-\frac{1}{2} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)^T \boldsymbol{\Sigma}_n^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_n)\right]$$

$$\boldsymbol{\Sigma}_n^{-1} = n\boldsymbol{\Sigma}^{-1} + \boldsymbol{\Sigma}_p^{-1}$$

$$\boldsymbol{\mu}_n = \boldsymbol{\Sigma}_p \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) + \frac{1}{n} \boldsymbol{\Sigma} \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \boldsymbol{\mu}_p$$

$$\boldsymbol{\Sigma}_n = \boldsymbol{\Sigma}_p \left( \boldsymbol{\Sigma}_p + \frac{1}{n} \boldsymbol{\Sigma} \right)^{-1} \frac{1}{n} \boldsymbol{\Sigma}$$

CS 2750 Machine Learning

## Other distributions

### Gamma distribution:

$$p(x | a, b) = \frac{1}{\Gamma(a) b^a} x^{a-1} e^{-\frac{x}{b}} \quad \text{for } x \in [0, \infty]$$

### Exponential distribution:

- A special case of Gamma for  $a=1$

$$p(x | b) = \left( \frac{1}{b} \right) e^{-\frac{x}{b}}$$

### Poisson distribution:

$$p(x | \lambda) = \frac{e^{-\lambda} \lambda^x}{x!} \quad \text{for } x \in \{0, 1, 2, \dots\}$$

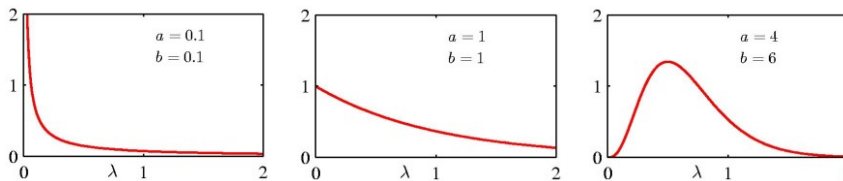
CS 2750 Machine Learning

## Other distributions

### Gamma distribution:

$$p(\lambda | a, b) = \frac{1}{\Gamma(a)b^a} \lambda^{a-1} e^{-\frac{\lambda}{b}} \quad \text{for } \lambda \in [0, \infty]$$

$$E(\lambda) = \frac{a}{b} \quad \text{var}(\lambda) = \frac{a}{b^2}$$



## Sequential Bayesian parameter estimation

- **Sequential Bayesian approach**

- Under the iid the estimates of the posterior can be computed incrementally for a sequence of data points

$$p(\Theta | D, \xi) = \frac{p(D | \Theta, \xi) p(\Theta | \xi)}{\int_{\Theta} p(D | \Theta, \xi) p(\Theta | \xi) d\Theta}$$

- If we use a conjugate prior we get back the same posterior
- Assume we split the data D in the last element  $\mathbf{x}$  and the rest  $p(D | \Theta) = P(\mathbf{x} | \Theta) P(D_{n-1} | \Theta)$

- **Then:** A “new” prior

$$p(\Theta | D, \xi) = \frac{P(\mathbf{x} | \Theta) \overbrace{P(D_{n-1} | \Theta) p(\Theta | \xi)}^{\text{A “new” prior}}}{\int_{\Theta} P(\mathbf{x} | \Theta) P(D_{n-1} | \Theta) p(\Theta | \xi) d\Theta}$$