

Parametric models

Maximum Likelihood and Bayesian Density Estimate

Muhammad Sarim

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Introduction

- Bayesian Decision Theory shows us how to design an optimal classifier if we know the prior probabilities $P(w_i)$ and the class-conditional densities $p(x|w_i)$.

Introduction

- Bayesian Decision Theory shows us how to design an optimal classifier if we know the prior probabilities $P(w_i)$ and the class-conditional densities $p(x|w_i)$.
- Unfortunately, we rarely have complete knowledge of the probabilistic structure.

Introduction

- Bayesian Decision Theory shows us how to design an optimal classifier if we know the prior probabilities $P(w_i)$ and the class-conditional densities $p(x|w_i)$.
- Unfortunately, we rarely have complete knowledge of the probabilistic structure.
- However, we can often find design samples or *training data* that include particular representatives of the patterns we want to classify.

Introduction

- To simplify the problem, we can parameterize the conditional densities and estimate these parameters using training data.

Introduction

- To simplify the problem, we can parameterize the conditional densities and estimate these parameters using training data.
- Then, we can use the resulting estimates as if they were the true values and perform classification using the Bayesian decision rule.

Introduction

- To simplify the problem, we can parameterize the conditional densities and estimate these parameters using training data.
- Then, we can use the resulting estimates as if they were the true values and perform classification using the Bayesian decision rule.
- We will consider only the supervised learning case where the true class label for each sample is known.

Introduction

- We will study two estimation procedures:

Introduction

- We will study two estimation procedures:
 - *Maximum likelihood estimation*
 - Views the parameters as quantities whose values are fixed but unknown
 - Estimate these values by maximizing the probability of obtaining the samples observed

Introduction

- We will study two estimation procedures:
 - *Maximum likelihood estimation*
 - Views the parameters as quantities whose values are fixed but unknown
 - Estimate these values by maximizing the probability of obtaining the samples observed
 - *Bayesian estimation*
 - Views the parameters as random variables having some known prior distribution
 - Observing new samples converts the prior to a posterior density

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Maximum Likelihood Estimation

- Suppose we have a set $D = \{x_1, \dots, x_n\}$ of independent and identically distributed (*i.i.d.*) samples drawn from the density $p(x|\theta)$.

Maximum Likelihood Estimation

- Suppose we have a set $D = \{x_1, \dots, x_n\}$ of independent and identically distributed (*i.i.d.*) samples drawn from the density $p(x|\theta)$.
- We would like to use training samples in D to estimate the unknown parameter vector θ .

Maximum Likelihood Estimation

- Suppose we have a set $D = \{x_1, \dots, x_n\}$ of independent and identically distributed (*i.i.d.*) samples drawn from the density $p(x|\theta)$.
- We would like to use training samples in D to estimate the unknown parameter vector θ .
- Define $L(\theta|D)$ as the *likelihood function* of θ with respect to D as

$$L(\theta|D) = p(D|\theta) = p(x_1, \dots, x_n|\theta) = \prod_{i=1}^n p(x_i|\theta)$$

Maximum Likelihood Estimation

- The *maximum likelihood estimate* (MLE) of θ is, by definition, the value $\hat{\theta}$ that maximizes $L(\theta|D)$ and can be computed as

$$\hat{\theta} = \arg \max_{\theta} L(\theta|D)$$

Maximum Likelihood Estimation

- The *maximum likelihood estimate* (MLE) of θ is, by definition, the value $\hat{\theta}$ that maximizes $L(\theta|D)$ and can be computed as

$$\hat{\theta} = \arg \max_{\theta} L(\theta|D)$$

- It is often easier to work with the logarithm of the likelihood function (*log-likelihood function*) that gives

$$\hat{\theta} = \arg \max_{\theta} \log L(\theta|D) = \arg \max_{\theta} \sum_{i=1}^n \log p(x_i|\theta)$$

Maximum Likelihood Estimation

- If the number of parameters is p , i.e.,
 $\theta = (\theta_1, \dots, \theta_p)^T$, define the gradient operator

$$\nabla_{\theta} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

Maximum Likelihood Estimation

- If the number of parameters is p , i.e., $\theta = (\theta_1, \dots, \theta_p)^T$, define the gradient operator

$$\nabla_{\theta} \equiv \begin{bmatrix} \frac{\partial}{\partial \theta_1} \\ \vdots \\ \frac{\partial}{\partial \theta_p} \end{bmatrix}$$

- Then, the MLE of θ should satisfy the necessary conditions

$$\nabla_{\theta} \log L(\theta|D) = \sum_{i=1}^n \nabla_{\theta} \log p(x_i|\theta) = 0$$

Maximum Likelihood Estimation

- Properties of MLEs:
 - The MLE is the parameter point for which the observed sample is the most likely.
 - The procedure with partial derivatives may result in several local extrema. We should check each solution individually to identify the global optimum.
 - Boundary conditions must also be checked separately for extrema.
 - Invariance property: if $\hat{\theta}$ is the MLE of θ , then for any function $f(\theta)$, the MLE of $f(\theta)$ is $f(\hat{\theta})$.

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

The Gaussian Case

- Suppose that $p(x|\theta) = N(\mu, \Sigma)$.

The Gaussian Case

- Suppose that $p(x|\theta) = N(\mu, \Sigma)$.
 - When Σ is known but μ is unknown:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

The Gaussian Case

- Suppose that $p(x|\theta) = N(\mu, \Sigma)$.
 - When Σ is known but μ is unknown:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

- When both μ and Σ are unknown:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and} \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

The Bernoulli Case

- Suppose that $P(x|\theta) = \text{Bernoulli}(\theta) = \theta^x(1 - \theta)^{1-x}$ where $x = 0, 1$ and $0 \leq \theta \leq 1$.

The Bernoulli Case

- Suppose that $P(x|\theta) = \text{Bernoulli}(\theta) = \theta^x(1 - \theta)^{1-x}$ where $x = 0, 1$ and $0 \leq \theta \leq 1$.
- The MLE of θ can be computed as

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n x_i$$

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Bias of Estimators

- *Bias* of an estimator $\hat{\theta}$ is the difference between the expected value of $\hat{\theta}$ and θ .

Bias of Estimators

- *Bias* of an estimator $\hat{\theta}$ is the difference between the expected value of $\hat{\theta}$ and θ .
- The MLE of μ is an unbiased estimator for μ because $E[\hat{\mu}] = \mu$.

Bias of Estimators

- *Bias* of an estimator $\hat{\theta}$ is the difference between the expected value of $\hat{\theta}$ and θ .
- The MLE of μ is an unbiased estimator for μ because $E[\hat{\mu}] = \mu$.
- The MLE of Σ is not an unbiased estimator for Σ because $E[\hat{\Sigma}] = \frac{n-1}{n}\Sigma \neq \Sigma$.

Bias of Estimators

- *Bias* of an estimator $\hat{\theta}$ is the difference between the expected value of $\hat{\theta}$ and θ .
- The MLE of μ is an unbiased estimator for μ because $E[\hat{\mu}] = \mu$.
- The MLE of Σ is not an unbiased estimator for Σ because $E[\hat{\Sigma}] = \frac{n-1}{n}\Sigma \neq \Sigma$.
- The *sample covariance*

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

is an unbiased estimator for Σ .

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Goodness-of-fit

- To measure how well a fitted distribution resembles the sample data (*goodness-of-fit*), we can use the Kolmogorov-Smirnov test statistic.

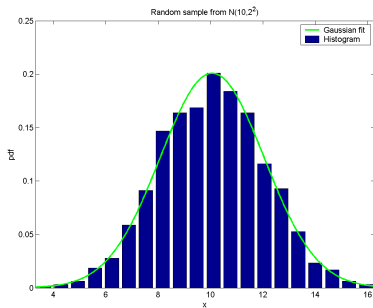
Goodness-of-fit

- To measure how well a fitted distribution resembles the sample data (*goodness-of-fit*), we can use the Kolmogorov-Smirnov test statistic.
- It is defined as the maximum value of the absolute difference between the cumulative distribution function estimated from the sample and the one calculated from the fitted distribution.

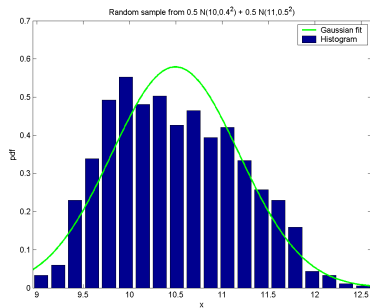
Goodness-of-fit

- To measure how well a fitted distribution resembles the sample data (*goodness-of-fit*), we can use the Kolmogorov-Smirnov test statistic.
- It is defined as the maximum value of the absolute difference between the cumulative distribution function estimated from the sample and the one calculated from the fitted distribution.
- After estimating the parameters for different distributions, we can compute the Kolmogorov-Smirnov statistic for each distribution and choose the one with the smallest value as the best fit to our sample.

Maximum Likelihood Estimation Examples

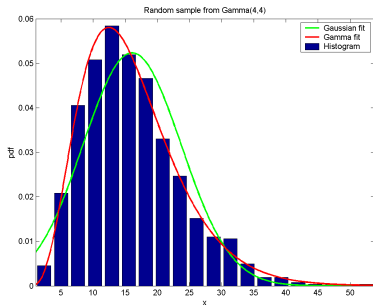


True pdf is $N(10, 4)$. Estimated pdf is $N(10.1, 3.9)$.

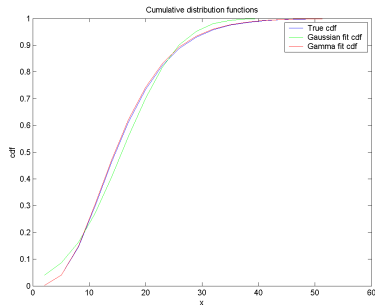


True pdf is $0.5N(10, 0.16) + 0.5N(11, 0.25)$. Estimated pdf is $N(10.5, 0.5)$.

Maximum Likelihood Estimation Examples



True pdf is Gamma(4,4). Estimated pdfs are $N(15.8, 62.1)$ and Gamma(4.0, 3.9).



Cumulative distribution functions.

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Bayesian Estimation

- Assume that θ is a quantity whose variation can be described by the prior probability distribution $p(\theta)$.

Bayesian Estimation

- Assume that θ is a quantity whose variation can be described by the prior probability distribution $p(\theta)$.
- Suppose the set $D = \{x_1, \dots, x_n\}$ contains the samples drawn independently from the density $p(x|\theta)$ whose form is assumed to be known but θ is not known exactly.

Bayesian Estimation

- Given D , the prior distribution can be updated to form the posterior distribution using the Bayes rule

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

where

$$p(D) = \int p(D|\theta) p(\theta) d\theta$$

and

$$p(D|\theta) = \prod_{i=1}^n p(x_i|\theta)$$

Bayesian Estimation

- The posterior distribution $p(\theta|D)$ can be used to find estimates for θ (e.g., the expected value of $p(\theta|D)$ can be used as an estimate for θ).

Bayesian Estimation

- The posterior distribution $p(\theta|D)$ can be used to find estimates for θ (e.g., the expected value of $p(\theta|D)$ can be used as an estimate for θ).
- Then, the conditional density $p(x|D)$ can be computed as

$$p(x|D) = \int p(x|\theta) p(\theta|D) d\theta$$

and can be used in the Bayesian classifier.

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

The Gaussian Case

- Consider the univariate case $p(x|\mu) = N(\mu, \sigma^2)$ where μ is the only unknown parameter with a prior distribution $p(\mu) = N(\mu_0, \sigma_0^2)$ (σ^2 , μ_0 and σ_0^2 are all known).

The Gaussian Case

- Consider the univariate case $p(x|\mu) = N(\mu, \sigma^2)$ where μ is the only unknown parameter with a prior distribution $p(\mu) = N(\mu_0, \sigma_0^2)$ (σ^2 , μ_0 and σ_0^2 are all known).
- This corresponds to drawing a value for μ from the population with density $p(\mu)$, treating it as the true value in the density $p(x|\mu)$, and drawing samples for x from this density.

The Gaussian Case

- Given $D = \{x_1, \dots, x_n\}$, we obtain

$$\begin{aligned} p(\mu|D) &\propto \prod_{i=1}^n p(x_i|\mu)p(\mu) \\ &\propto \exp \left[-\frac{1}{2} \left(\left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2} \right) \mu^2 - 2 \left(\frac{1}{\sigma^2} \sum_{i=1}^n x_i + \frac{\mu_0}{\sigma_0^2} \right) \mu \right) \right] \\ &= N(\mu_n, \sigma_n^2) \end{aligned}$$

where

$$\begin{aligned} \mu_n &= \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2} \right) \hat{\mu}_n + \left(\frac{\sigma^2}{n\sigma_0^2 + \sigma^2} \right) \mu_0 & \left(\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i \right) \\ \sigma_n^2 &= \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2} \end{aligned}$$

The Gaussian Case

- μ_0 is our best prior guess and σ_0^2 is the uncertainty about this guess.

The Gaussian Case

- μ_0 is our best prior guess and σ_0^2 is the uncertainty about this guess.
- μ_n is our best guess after observing D and σ_n^2 is the uncertainty about this guess.

The Gaussian Case

- μ_0 is our best prior guess and σ_0^2 is the uncertainty about this guess.
- μ_n is our best guess after observing D and σ_n^2 is the uncertainty about this guess.
- μ_n always lies between $\hat{\mu}_n$ and μ_0 .
 - If $\sigma_0 = 0$, then $\mu_n = \mu_0$ (no observation can change our prior opinion).
 - If $\sigma_0 \gg \sigma$, then $\mu_n = \hat{\mu}_n$ (we are very uncertain about our prior guess).
 - Otherwise, μ_n approaches $\hat{\mu}_n$ as n approaches infinity.

The Gaussian Case

- Given the posterior density $p(\mu|D)$, the conditional density $p(x|D)$ can be computed as

$$p(x|D) = N(\mu_n, \sigma^2 + \sigma_n^2)$$

where the conditional mean μ_n is treated as if it were the true mean, and the known variance is increased to account for our lack of exact knowledge of the mean μ .

The Gaussian Case

- Consider the multivariate case $p(x|\mu) = N(\mu, \Sigma)$ where μ is the only unknown parameter with a prior distribution $p(\mu) = N(\mu_0, \Sigma_0)$ (Σ , μ_0 and Σ_0 are all known).

The Gaussian Case

- Consider the multivariate case $p(x|\mu) = N(\mu, \Sigma)$ where μ is the only unknown parameter with a prior distribution $p(\mu) = N(\mu_0, \Sigma_0)$ (Σ , μ_0 and Σ_0 are all known).
- Given $D = \{x_1, \dots, x_n\}$, we obtain

$$p(\mu|D) \propto \exp \left[-\frac{1}{2} \left(\mu^T \left(n\Sigma^{-1} + \Sigma_0^{-1} \right) \mu - 2\mu^T \left(\Sigma^{-1} \sum_{i=1}^n x_i + \Sigma_0^{-1} \mu_0 \right) \right) \right]$$

The Gaussian Case

- It follows that

$$p(\mu|D) = N(\mu_n, \Sigma_n)$$

where

$$\mu_n = \Sigma_0 \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \hat{\mu}_n + \frac{1}{n} \Sigma \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \mu_0$$
$$\Sigma_n = \frac{1}{n} \Sigma_0 \left(\Sigma_0 + \frac{1}{n} \Sigma \right)^{-1} \Sigma$$

The Gaussian Case

- Given the posterior density $p(\mu|D)$, the conditional density $p(x|D)$ can be computed as

$$p(x|D) = N(\mu_n, \Sigma + \Sigma_n)$$

which can be viewed as the sum of a random vector μ with $p(\mu|D) = N(\mu_n, \Sigma_n)$ and an independent random vector y with $p(y) = N(0, \Sigma)$.

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

The Bernoulli Case

- Consider $P(x|\theta) = \text{Bernoulli}(\theta)$ where θ is the unknown parameter with a prior distribution $p(\theta) = \text{Beta}(\alpha, \beta)$ (α and β are both known).

The Bernoulli Case

- Consider $P(x|\theta) = \text{Bernoulli}(\theta)$ where θ is the unknown parameter with a prior distribution $p(\theta) = \text{Beta}(\alpha, \beta)$ (α and β are both known).
- Given $D = \{x_1, \dots, x_n\}$, we obtain

$$p(\theta|D) = \text{Beta} \left(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i \right)$$

The Bernoulli Case

- The Bayes estimate of θ can be computed as the expected value of $p(\theta|D)$

$$\begin{aligned}\hat{\theta} &= \frac{\alpha + \sum_{i=1}^n x_i}{\alpha + \beta + n} \\ &= \left(\frac{n}{\alpha + \beta + n} \right) \frac{1}{n} \sum_{i=1}^n x_i + \left(\frac{\alpha + \beta}{\alpha + \beta + n} \right) \frac{\alpha}{\alpha + \beta}\end{aligned}$$

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Conjugate Priors

- A *conjugate prior* is one which, when multiplied with the probability of the observation, gives a posterior probability having the same functional form as the prior.

Conjugate Priors

- A *conjugate prior* is one which, when multiplied with the probability of the observation, gives a posterior probability having the same functional form as the prior.
- This relationship allows the posterior to be used as a prior in further computations.

Conjugate Priors

- A *conjugate prior* is one which, when multiplied with the probability of the observation, gives a posterior probability having the same functional form as the prior.
- This relationship allows the posterior to be used as a prior in further computations.

Table: Conjugate prior distributions.

<i>pdf generating the sample</i>	<i>corresponding conjugate prior</i>
Normal	Normal
Exponential	Gamma
Poisson	Gamma
Binomial	Beta
Multinomial	Dirichlet

Contents

1 Introduction

2 MLE

- Examples
- Bias
- GoF

3 Bayesian Estimation

- The Gaussian Case
- The Bernoulli Case
- Conjugate Priors
- Recursive Bayes Learning

Recursive Bayes Learning

- What about the convergence of $p(x|D)$ to $p(x)$?

Recursive Bayes Learning

- What about the convergence of $p(x|D)$ to $p(x)$?
- Given $D^n = \{x_1, \dots, x_n\}$, for $n > 1$

$$p(D^n|\theta) = p(x_n|\theta)p(D^{n-1}|\theta)$$

and

$$p(\theta|D^n) = \frac{p(x_n|\theta) p(\theta|D^{n-1})}{\int p(x_n|\theta) p(\theta|D^{n-1}) d\theta}$$

where

$$p(\theta|D^0) = p(\theta)$$

Recursive Bayes Learning

- What about the convergence of $p(x|D)$ to $p(x)$?
- Given $D^n = \{x_1, \dots, x_n\}$, for $n > 1$

$$p(D^n|\theta) = p(x_n|\theta)p(D^{n-1}|\theta)$$

and

$$p(\theta|D^n) = \frac{p(x_n|\theta) p(\theta|D^{n-1})}{\int p(x_n|\theta) p(\theta|D^{n-1}) d\theta}$$

where

$$p(\theta|D^0) = p(\theta)$$

⇒ quite useful if the distributions can be represented using only a few parameters (*sufficient statistics*)