

CS 3750 Machine Learning

Lecture 2

Density estimation

Milos Hauskrecht

milos@cs.pitt.edu

Sennott Square, x4-8845

<http://www.cs.pitt.edu/~milos/courses/cs3750/>

CS 3750 Advanced Machine Learning

Density estimation

Data: $D = \{D_1, D_2, \dots, D_n\}$

$D_i = \mathbf{x}_i$ a vector of attribute values

Attributes:

- modeled by random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$ with:
 - **Continuous values**
 - **Discrete values**

E.g. *blood pressure* with numerical values

or *chest pain* with discrete values

[no-pain, mild, moderate, strong]

Underlying true probability distribution:

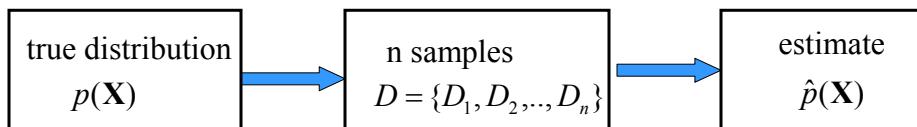
$p(\mathbf{X})$

CS 3750 Advanced Machine Learning

Density estimation

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

Objective: try to estimate the underlying true probability distribution over variables \mathbf{X} , $p(\mathbf{X})$, using examples in D



Standard (iid) assumptions: Samples

- are **independent** of each other
- come from the same (**identical**) **distribution** (fixed $p(\mathbf{X})$)

CS 3750 Advanced Machine Learning

Density estimation

Types of density estimation:

Parametric

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

Non-parametric

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

Semi-parametric

CS 3750 Advanced Machine Learning

Parametric density estimation

Parametric density estimation

Basic settings:

- A set of random variables $\mathbf{X} = \{X_1, X_2, \dots, X_d\}$
- **A model of the distribution** over variables in \mathbf{X} with parameters Θ
- **Data** $D = \{D_1, D_2, \dots, D_n\}$

Objective: find parameters $\hat{\Theta}$ that describe $p(\mathbf{X} | \Theta)$ the best

CS 3750 Advanced Machine Learning

Parameter learning

What is the best set of parameters?

- **Maximum likelihood (ML) estimates**

maximize $p(D | \Theta, \xi)$

ξ - represents prior (background) knowledge

- **Maximum a posteriori probability (MAP) estimate**

maximize $p(\Theta | D, \xi)$

Selects the mode of the posterior

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

CS 3750 Advanced Machine Learning

Parameter learning

- **Both ML or MAP pick one parameter value**
 - Is it always the best solution?
 - **Bayesian approach**
 - Remedies the limitation of one choice
 - Keeps and uses complete posterior distribution $p(\Theta | D, \xi)$
 - Optimization is replaced with integration
 - **How is it used? Assume we want:** $P(\mathbf{x} | D, \xi)$
 - Consider all parameter settings and averages the result

– **Example:** predict the result of the outcome $x=1$

$$P(x=1 \mid D, \xi)$$

CS 3750 Advanced Machine Learning

Bernoulli distribution.

Outcomes: x_i with values 0 or 1 (head or tail)

Data: D a sequence of outcomes x_i

Model: probability of an outcome 1 θ
probability of 0 $(1-\theta)$

$$P(x_i | \theta) = \theta^{x_i} (1-\theta)^{1-x_i} \quad \text{Bernoulli distribution}$$

CS 3750 Advanced Machine Learning

Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D | \theta, \xi) = \prod_{i=1}^n \theta^{x_i} (1-\theta)^{(1-x_i)}$$

Maximum likelihood estimate

$$\theta_{ML} = \arg \max_{\theta} P(D | \theta, \xi)$$

Optimize log-likelihood

$$l(D, \theta) = \log P(D | \theta, \xi) = \log \prod_{i=1}^n \theta^{x_i} (1-\theta)^{(1-x_i)} =$$
$$\sum_{i=1}^n x_i \log \theta + (1-x_i) \log (1-\theta) = \underbrace{\log \theta \sum_{i=1}^n x_i}_{N_1 - \text{number of 1s seen}} + \underbrace{\log (1-\theta) \sum_{i=1}^n (1-x_i)}_{N_2 - \text{number of 0s seen}}$$

CS 3750 Advanced Machine Learning

Maximum likelihood (ML) estimate.

Optimize log-likelihood

$$l(D, \theta) = N_1 \log \theta + N_2 \log (1-\theta)$$

Set derivative to zero

$$\frac{\partial l(D, \theta)}{\partial \theta} = \frac{N_1}{\theta} - \frac{N_2}{(1-\theta)} = 0$$

Solving

$$\theta = \frac{N_1}{N_1 + N_2}$$

ML Solution: $\theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$

CS 3750 Advanced Machine Learning

Maximum a posteriori estimate

Maximum a posteriori estimate

- Selects the mode of the posterior distribution

$$\theta_{MAP} = \arg \max_{\theta} p(\theta | D, \xi)$$

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) p(\theta | \xi)}{P(D | \xi)} \quad (\text{via Bayes rule})$$

$P(D | \theta, \xi)$ - is the likelihood of data

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

$p(\theta | \xi)$ - is the prior probability on θ

How to choose the prior probability?

CS 3750 Advanced Machine Learning

Prior distribution

Choice of prior: Beta distribution

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

Why?

Beta distribution “fits” binomial sampling - conjugate choices

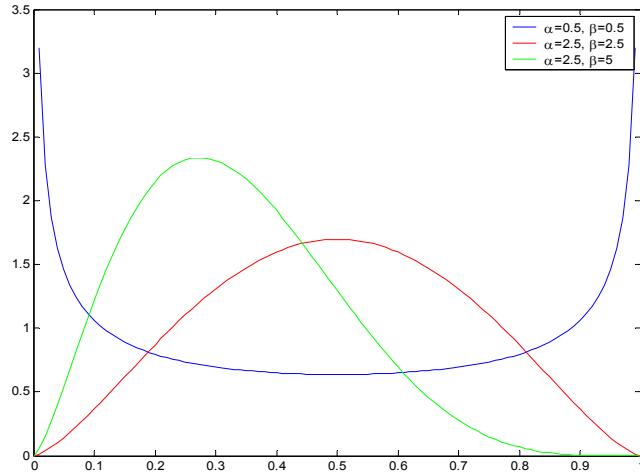
$$P(D | \theta, \xi) = \theta^{N_1} (1-\theta)^{N_2}$$

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

MAP Solution: $\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$

CS 3750 Advanced Machine Learning

Beta distribution



CS 3750 Advanced Machine Learning

Bayesian approach

- **Posterior probability:**

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

- **Probability of an outcome $x=1$ in the next trial**

$$\begin{aligned} P(x=1 | D, \xi) &= \int_0^1 P(x=1 | \theta, \xi) p(\theta | D, \xi) d\theta \\ &= \int_0^1 \theta p(\theta | D, \xi) d\theta = E(\theta) \end{aligned}$$

- **Equivalent to the expected value of the parameter**

- expectation is taken with regard to the posterior distribution

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

CS 3750 Advanced Machine Learning

Bayesian learning

Expected value of the parameter

$$\begin{aligned} E(\theta) &= \int_0^1 \theta \text{Beta}(\theta | \eta_1, \eta_2) d\theta = \int_0^1 \theta \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \theta^{\eta_1-1} (1-\theta)^{\eta_2-1} d\theta \\ &= \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \int_0^1 \theta^{\eta_1} (1-\theta)^{\eta_2-1} d\theta \\ &= \frac{\Gamma(\eta_1 + \eta_2)}{\Gamma(\eta_1)\Gamma(\eta_2)} \frac{\Gamma(\eta_1+1)\Gamma(\eta_2)}{\Gamma(\eta_1+\eta_2+1)} \underbrace{\int_0^1 \text{Beta}(\eta_1+1, \eta_2) d\theta}_1 \\ &= \frac{\eta_1}{\eta_1 + \eta_2} \end{aligned}$$

Note: $\Gamma(\alpha + 1) = \alpha\Gamma(\alpha)$

CS 3750 Advanced Machine Learning

Expected value

- Predictive probability of an outcome $x=1$ in the next trial

$$P(x=1 | D, \xi) = E(\theta)$$

- Substituting the results for

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

- We get

$$P(x=1 | D, \xi) = E(\theta) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$$

- Instead of MAP and ML choice of the parameter we can use the expected value of the parameter

$$\hat{\theta} = E(\theta)$$

CS 3750 Advanced Machine Learning

Binomial distribution.

Data: D – a set of order-independent outcomes with two possible values – 0 or 1 (head or tail)

N_1 - number of heads seen N_2 - number of tails seen

We treat D as a multi-set !!!

Model: probability of a 1 θ
probability of a 0 $(1-\theta)$

Probability of an outcome

$$P(D | \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} \quad \text{Binomial distribution}$$

CS 3750 Advanced Machine Learning

Maximum likelihood (ML) estimate.

Likelihood of data:

$$P(D | \theta) = \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} = \frac{N!}{N_1! N_2!} \theta^{N_1} (1-\theta)^{N_2}$$

Log-likelihood

$$l(D, \theta) = \log \binom{N}{N_1} \theta^{N_1} (1-\theta)^{N_2} = \log \frac{N!}{N_1! N_2!} + N_1 \log \theta + N_2 \log (1-\theta)$$

Constant from the point of optimization !!!

$$\text{ML Solution: } \theta_{ML} = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$$

The same as for Bernoulli and D **with iid** sequence of examples

CS 3750 Advanced Machine Learning

Maximum a posteriori estimate

MAP estimate $\theta_{MAP} = \arg \max_{\theta} p(\theta | D, \xi)$

$$p(\theta | D, \xi) = \frac{P(D | \theta, \xi) p(\theta | \xi)}{P(D | \xi)} \quad (\text{via Bayes rule})$$

Prior choice

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

Likelihood

$$P(D | \theta) = \frac{\Gamma(N_1 + N_2)}{\Gamma(N_1)\Gamma(N_2)} \theta^{N_1} (1-\theta)^{N_2}$$

Posterior

$$p(\theta | D, \xi) = Beta(\alpha_1 + N_1, \alpha_2 + N_2)$$

MAP estimate

$$\theta_{MAP} = \frac{\alpha_1 + N_1 - 1}{\alpha_1 + \alpha_2 + N_1 + N_2 - 2}$$

CS 3750 Advanced Machine Learning

Bayesian learning, expectation

The result is the same as for Bernoulli distribution

$$E(\theta) = \int_0^1 \theta Beta(\theta | \eta_1, \eta_2) d\theta = \frac{\eta_1}{\eta_1 + \eta_2}$$

Expected value of the parameter

$$E(\theta) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$$

Predictive probability of an event $x=1$

$$E(\theta) = P(x=1 | \theta, \xi) = \frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \alpha_2 + N_2}$$

CS 3750 Advanced Machine Learning

Multinomial distribution

A kind of multi-way coin toss (roll of dice)

- **Data:** a set of N outcomes (multi-set)

N_i - a number of times an outcome i has been seen

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

Probability of data (likelihood)

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

Multinomial distribution

ML estimate:

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

CS 3750 Advanced Machine Learning

MAP estimate

Choice of 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 multinomial

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

Posterior distribution

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

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

CS 3750 Advanced Machine Learning

Bayesian learning

The result is analogous to the result for binomial

$$E(\boldsymbol{\theta}) = \iint_{\substack{0 \leq \theta_i \leq 1, \\ \sum \theta_i = 1}} \boldsymbol{\theta} \text{Dir}(\boldsymbol{\theta} | \boldsymbol{\eta}) d\boldsymbol{\theta} = \left(\frac{\eta_1}{\eta_1 + \eta_2 + \eta_k}, \dots, \frac{\eta_i}{\eta_1 + \eta_2 + \eta_k}, \dots, \frac{\eta_k}{\eta_1 + \eta_2 + \eta_k} \right)$$

Bayesian estimate substitutes posterior

$$E(\boldsymbol{\theta}) = \left(\frac{\alpha_1 + N_1}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}, \dots, \frac{\alpha_i + N_i}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}, \dots, \frac{\alpha_k + N_k}{\alpha_1 + N_1 + \dots + \alpha_k + N_k} \right)$$

Represents the predictive probability of an event $x=i$

$$P(x=i | \boldsymbol{\theta}, \xi) = \frac{\alpha_i + N_i}{\alpha_1 + N_1 + \dots + \alpha_k + N_k}$$

CS 3750 Advanced Machine Learning

Other distributions

The same ideas can be applied to other distributions

- Typically we choose distributions that behave in a nice way so that the 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 - Gamma
- Gaussian - Gaussian (mean) and Wishart (covariance)

CS 3750 Advanced Machine Learning

Distributions from the exponential family

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:

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

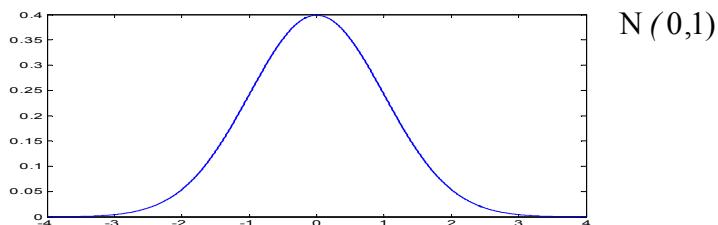
CS 3750 Advanced Machine Learning

Gaussian (normal) distribution

- **Gaussian:** $x \sim N(\mu, \sigma)$
- **Parameters:** μ - mean
 σ - 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:**



CS 3750 Advanced Machine Learning

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

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

- **Unbiased estimate:**

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