

## CS 2750 Machine Learning Lecture 5

### Density estimation

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

---

CS 2750 Machine Learning

### Density estimations

#### Topics:

- **Density estimation:** ✓
    - Maximum likelihood (ML)
    - Bayesian parameter estimates
    - MAP
  - **Bernoulli distribution.** ✓
  - **Binomial distribution** ✓
  - **Multinomial distribution** ✓
  - **Normal distribution** ✓
  - **Exponential family**
- Nonparametric family**

---

CS 2750 Machine Learning

## Parametric density estimation

### Parametric density estimation:

- 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**  $\Theta : \hat{p}(\mathbf{X} | \Theta)$
- **Data**  $D = \{D_1, D_2, \dots, D_n\}$

### Objective:

find parameters  $\Theta$  such that  $p(\mathbf{X} | \Theta)$  describes data  $D$  the best

---

CS 2750 Machine Learning

## Parameter estimation (learning)

- **Maximum likelihood (ML)**  
 $\Theta_{ML} = \arg \max_{\Theta} p(D | \Theta, \xi)$
- **Maximum a posteriori probability (MAP)**  
 $\Theta_{MAP} = \arg \max_{\Theta} p(\Theta | D, \xi)$
- **Bayesian parameter estimation**
  - use the posterior density  
 $p(\Theta | D, \xi)$
- **Expected value**

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

---

CS 2750 Machine Learning

## Exponential family of distribution

### Exponential family of distributions

- well behaved distributions with respect to ML and Bayesian updating

**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)**

CS 2750 Machine Learning

## 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  $x$  and the rest  $p(D | \Theta) = P(x | \Theta) P(D_{n-1} | \Theta)$

### • Then:

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

CS 2750 Machine Learning

## Exponential family

### Exponential family:

- all probability mass / density functions that can be written in the exponential normal form

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- $\boldsymbol{\eta}$  a vector of **natural (or canonical) parameters**
- $t(\mathbf{x})$  a function referred to as a **sufficient statistic**
- $h(\mathbf{x})$  a function of  $\mathbf{x}$  (it is less important)
- $Z(\boldsymbol{\eta})$  a normalization constant (a **partition function**)

$$Z(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}$$

- Other common form:

$$f(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta})] \quad \log Z(\boldsymbol{\eta}) = A(\boldsymbol{\eta})$$

CS 2750 Machine Learning

## Exponential family: examples

- Bernoulli distribution**

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp\left\{ \log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi) \right\} \\ &= \exp\{\log(1-\pi)\} \exp\left\{ \log\left(\frac{\pi}{1-\pi}\right)x \right\} \end{aligned}$$

- Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- Parameters**

$$\boldsymbol{\eta} = ?$$

$$t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ?$$

$$h(\mathbf{x}) = ?$$

CS 2750 Machine Learning

## Exponential family: examples

- **Bernoulli distribution**

$$\begin{aligned} p(x | \pi) &= \pi^x (1 - \pi)^{1-x} \\ &= \exp\left\{\log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi)\right\} \\ &= \exp\{\log(1-\pi)\} \exp\left\{\log\left(\frac{\pi}{1-\pi}\right)x\right\} \end{aligned}$$

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$$

- **Parameters**

$$\begin{aligned} \boldsymbol{\eta} &= \log \frac{\pi}{1-\pi} & t(\mathbf{x}) &= x \\ Z(\boldsymbol{\eta}) &= \frac{1}{1-\pi} = 1 + e^\eta & h(\mathbf{x}) &= 1 \end{aligned}$$

CS 2750 Machine Learning

## Exponential family: examples

- **Univariate Gaussian distribution**

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

- **Exponential family**

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(x) \exp[\boldsymbol{\eta}^T t(x)]$$

- **Parameters**

$$\begin{aligned} \boldsymbol{\eta} &= ? & t(\mathbf{x}) &= ? \\ Z(\boldsymbol{\eta}) &= ? & h(\mathbf{x}) &= ? \end{aligned}$$

CS 2750 Machine Learning

## Exponential family: examples

- **Univariate Gaussian distribution**

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

$$= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\mu}{2\sigma^2} - \log \sigma\right) \exp\left\{\frac{\mu}{\sigma^2} x - \frac{1}{2\sigma^2} x^2\right\}$$

- **Exponential family**  $f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T t(\mathbf{x})]$

- **Parameters**

$$\boldsymbol{\eta} = \begin{bmatrix} \mu / 2\sigma^2 \\ -1 / 2\sigma^2 \end{bmatrix} \quad t(\mathbf{x}) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$$

$$Z(\boldsymbol{\eta}) = \exp\left\{\frac{\mu}{2\sigma^2} + \log \sigma\right\} = \exp\left\{-\frac{\eta_1^2}{4\eta_2} - \frac{1}{2} \log(-2\eta_2)\right\}$$

$$h(\mathbf{x}) = 1 / \sqrt{2\pi}$$

CS 2750 Machine Learning

## Exponential family

- **For iid samples, the likelihood of data is**

$$P(D | \boldsymbol{\eta}) = \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\eta}) = \prod_{i=1}^n h(\mathbf{x}_i) \exp[\boldsymbol{\eta}^T t(\mathbf{x}_i) - A(\boldsymbol{\eta})]$$

$$= \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[ \sum_{i=1}^n \boldsymbol{\eta}^T t(\mathbf{x}_i) - nA(\boldsymbol{\eta}) \right]$$

$$= \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp\left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right]$$

- **Important:**

- the dimensionality of the sufficient statistic remains the same for different sample sizes (that is, different number of examples in D)

CS 2750 Machine Learning

## Exponential family

- The log likelihood of data is

$$\begin{aligned}l(D, \boldsymbol{\eta}) &= \log \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] \exp \left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right] \\ &= \log \left[ \prod_{i=1}^n h(\mathbf{x}_i) \right] + \left[ \boldsymbol{\eta}^T \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - nA(\boldsymbol{\eta}) \right]\end{aligned}$$

- Optimizing the loglikelihood

$$\nabla_{\boldsymbol{\eta}} l(D, \boldsymbol{\eta}) = \left( \sum_{i=1}^n t(\mathbf{x}_i) \right) - n \nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \mathbf{0}$$

- For the ML estimate it must hold

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{1}{n} \left( \sum_{i=1}^n t(\mathbf{x}_i) \right)$$

CS 2750 Machine Learning

## Exponential family

- Rewriting the gradient:

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log Z(\boldsymbol{\eta}) = \nabla_{\boldsymbol{\eta}} \log \int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \frac{\int t(\mathbf{x}) h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}}{\int h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x})\} d\mathbf{x}}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = \int t(\mathbf{x}) h(\mathbf{x}) \exp\{\boldsymbol{\eta}^T t(\mathbf{x}) - A(\boldsymbol{\eta})\} d\mathbf{x}$$

$$\nabla_{\boldsymbol{\eta}} A(\boldsymbol{\eta}) = E(t(\mathbf{x}))$$

- Result: 
$$E(t(\mathbf{x})) = \frac{1}{n} \left( \sum_{i=1}^n t(\mathbf{x}_i) \right)$$

- For the ML estimate, the parameters  $\boldsymbol{\eta}$  should be adjusted such that the expectation of the statistic  $t(\mathbf{x})$  is equal to the observed sample statistics

CS 2750 Machine Learning

## Moments of the distribution

- **For the exponential family**

- The k-th moment of the statistic corresponds to the k-th derivative of  $A(\boldsymbol{\eta})$
- If  $x$  is a component of  $t(x)$  then we get the moments of the distribution by differentiating its corresponding natural parameter

- **Example: Bernoulli**  $p(x | \pi) = \exp\left\{\log\left(\frac{\pi}{1-\pi}\right)x + \log(1-\pi)\right\}$   
 $A(\boldsymbol{\eta}) = \log\frac{1}{1-\pi} = \log(1+e^\eta)$

- **Derivatives:**

$$\frac{\partial A(\boldsymbol{\eta})}{\partial \eta} = \frac{\partial}{\partial \eta} \log(1+e^\eta) = \frac{e^\eta}{(1+e^\eta)} = \frac{1}{(1+e^{-\eta})} = \pi$$
$$\frac{\partial A(\boldsymbol{\eta})}{\partial \eta^2} = \frac{\partial}{\partial \eta} \frac{1}{(1+e^{-\eta})} = \pi(1-\pi)$$

CS 2750 Machine Learning

## Exponential family of distribution

### Bayesian parameter estimate

We have seen 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)**

CS 2750 Machine Learning



## Conjugate priors

For any member of the exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x})]$$

there exists a prior:

$$p(\boldsymbol{\eta} | \boldsymbol{\chi}, \nu) = u(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp[\nu \boldsymbol{\eta}^T \boldsymbol{\chi}]$$

Such that for  $n$  examples, the posterior is

$$p(\boldsymbol{\eta} | D, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+n} \exp\left[\boldsymbol{\eta}^T \left(\left[\sum_{i=1}^n \mathbf{t}(x_i)\right] + \nu \boldsymbol{\chi}\right)\right]$$

Note that:

$$P(D | \boldsymbol{\eta}) = \left(\frac{1}{Z(\boldsymbol{\eta})}\right)^n \left[\prod_{i=1}^n h(\mathbf{x}_i)\right] \exp\left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n \mathbf{t}(\mathbf{x}_i)\right)\right]$$

CS 2750 Machine Learning

## Conjugate priors

For any member of the exponential family

$$f(\mathbf{x} | \boldsymbol{\eta}) = \frac{1}{Z(\boldsymbol{\eta})} h(\mathbf{x}) \exp[\boldsymbol{\eta}^T \mathbf{t}(\mathbf{x})]$$

there exists a prior:

$$p(\boldsymbol{\eta} | \boldsymbol{\chi}, \nu) = u(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp[\nu \boldsymbol{\eta}^T \boldsymbol{\chi}]$$

Such that for  $n$  examples, the posterior is

$$p(\boldsymbol{\eta} | D, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{\nu+n} \exp\left[\boldsymbol{\eta}^T \left(\left[\sum_{i=1}^n \mathbf{t}(x_i)\right] + \nu \boldsymbol{\chi}\right)\right]$$

Pseudo-observation

Note that:

$$P(D | \boldsymbol{\eta}) = \left(\frac{1}{Z(\boldsymbol{\eta})}\right)^n \left[\prod_{i=1}^n h(\mathbf{x}_i)\right] \exp\left[\boldsymbol{\eta}^T \left(\sum_{i=1}^n \mathbf{t}(\mathbf{x}_i)\right)\right]$$

CS 2750 Machine Learning

## Nonparametric Methods

- **Parametric distribution models** are:
  - restricted to specific forms, which may not always be suitable;
  - Example: modelling a multimodal distribution with a single, unimodal model.
- **Nonparametric approaches:**
  - make few assumptions about the overall shape of the distribution being modelled.

CS 2750 Machine Learning

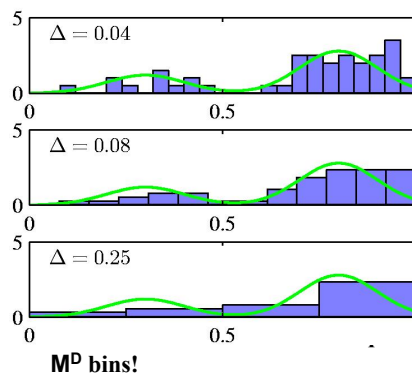
## Nonparametric Methods

### Histogram methods:

partition the data space into distinct bins with widths  $\Delta_i$  and count the number of observations,  $n_i$ , in each bin.

$$p_i = \frac{n_i}{N\Delta_i}$$

- Often, the same width is used for all bins,  $\Delta_i = \Delta$ .
- $\Delta$  acts as a smoothing parameter.



CS 2750 Machine Learning

## Nonparametric Methods

- Assume observations drawn from a density  $p(x)$  and consider a small region  $R$  containing  $x$  such that

$$P = \int_R p(x) dx$$

- The probability that  $K$  out of  $N$  observations lie inside  $R$  is  $\text{Bin}(K, N, P)$  and if  $N$  is large

$$K \cong NP$$

If the volume of  $R$ ,  $V$ , is sufficiently small,  $p(x)$  is approximately constant over  $R$  and

$$P \cong p(x)V$$

Thus

$$p(x) = \frac{P}{V}$$

$$p(x) = \frac{K}{NV}$$

## Nonparametric Methods: kernel methods

### Kernel Density Estimation:

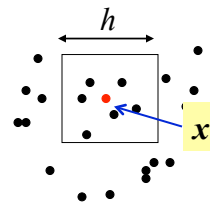
**Fix  $V$ , estimate  $K$  from the data.** Let  $R$  be a hypercube centred on  $\mathbf{x}$  and define the kernel function (Parzen window)

$$k\left(\frac{x - x_n}{h}\right) = \begin{cases} 1 & |(x_i - x_{ni})| / h \leq 1/2 \\ 0 & \text{otherwise} \end{cases} \quad i = 1, \dots, D$$

- It follows that

- and hence 
$$K = \sum_{n=1}^N k\left(\frac{x - x_n}{h}\right)$$

$$p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{x - x_n}{h}\right)$$



## Nonparametric Methods: smooth kernels

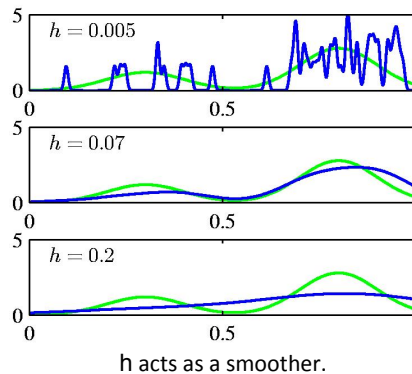
To avoid discontinuities in  $p(\mathbf{x})$  because of sharp boundaries use a **smooth kernel**, e.g. a Gaussian

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp\left\{-\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2}\right\}$$

- Any kernel such that

$$k(\mathbf{u}) \geq 0, \\ \int k(\mathbf{u}) d\mathbf{u} = 1$$

- will work.



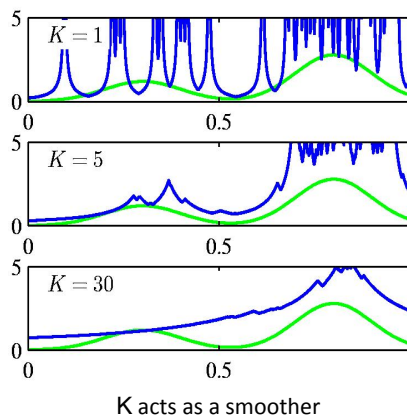
CS 2750 Machine Learning

## Nonparametric Methods: kNN estimation

### Nearest Neighbour Density Estimation:

fix  $K$ , estimate  $V$  from the data. Consider a hyper-sphere centred on  $\mathbf{x}$  and let it grow to a volume,  $V^*$ , that includes  $K$  of the given  $N$  data points. Then

$$p(\mathbf{x}) \simeq \frac{K}{NV^*}.$$



CS 2750 Machine Learning