

CS 2750 Machine Learning Lecture 7

Density estimation III

Milos Hauskrecht

milos@pitt.edu

5329 Sennott Square

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:

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

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**

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

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**

$$\begin{array}{ll}
 \boldsymbol{\eta} = \log \frac{\pi}{1 - \pi} \quad \leftarrow \text{logit function} & t(\mathbf{x}) = x \\
 Z(\boldsymbol{\eta}) = \frac{1}{1 - \pi} = 1 + e^\eta & h(\mathbf{x}) = 1
 \end{array}$$

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**

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

- **Parameters**

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

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

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

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

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**

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

- **Parameters**

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

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

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

Exponential family

- For iid samples, the likelihood of data is

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

- **Important:**
 - the dimensionality of the sufficient statistic remains the same with the number of samples

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

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

Exponential family

- Rewriting the gradient:

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

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:**

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

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

Non-parametric density estimation

Nonparametric Density Estimation

- **Parametric distribution models** are:
 - restricted to specific functional forms, which may not always be suitable;
 - **Example:** modeling a multimodal distribution with a single, unimodal model.



- **Nonparametric approaches:**
 - Do not make any strong assumption about the overall shape of the distribution being modelled.

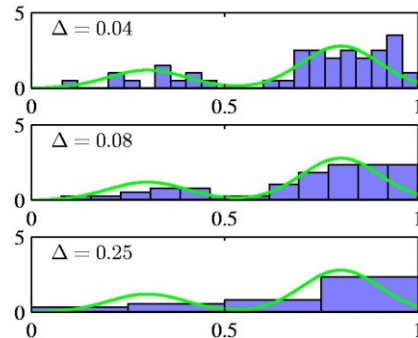
Nonparametric Methods

Histogram methods:

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

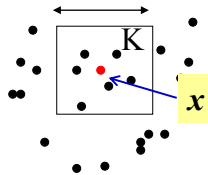
$$p_i = \frac{n_i}{N\Delta_i}$$

- Often, the same width is used for all bins, $\Delta_i = \Delta$.
- Δ acts as a smoothing parameter.
- Binning does not work well in the in a d -dimensional space,



Nonparametric Methods

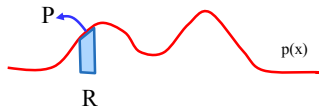
- Binning does not work well in the in a d-dimensional space,
 - M bins in each dimension will require M^d bins!
- **Solution:**
 - Build the estimates of $p(\mathbf{x})$ by considering the data points in D and how similar (or close) they are to \mathbf{x}
 - **Example: Parzen window**
 - As if we build a bin dynamically for \mathbf{x} for which we need $p(\mathbf{x})$



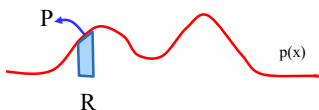
Nonparametric Methods

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

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

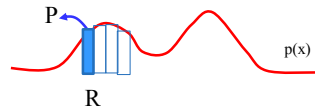


- The probability that K out of N observations lie inside R is $Bin(K, N, P)$ and if N is large $K \cong NP$



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

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



Thus

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

Putting things together we get:

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

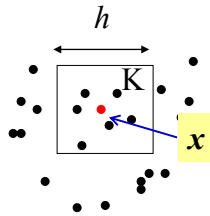
Nonparametric methods: kernel methods

Solution 1: Estimate the probability for \mathbf{x} based on the fixed volume V built around \mathbf{x}

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

- Fix V , estimate K from the data

Example: **Parzen window**



Nonparametric methods: kernel methods

Kernel Density Estimation:

- **Parzen window:** Let R be a hypercube centred on \mathbf{x} that defines the **kernel function:**

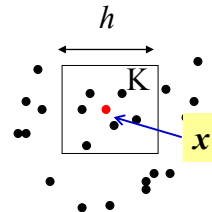
$$k\left(\frac{x - x_n}{h}\right) = \begin{cases} 1 & |x_i - x_{ni}| / h \leq 1/2 \\ 0 & \text{otherwise} \end{cases} \quad i = 1, \dots, D$$

- It follows that

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

- and hence

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



Nonparametric Methods: smooth kernels

To avoid discontinuities in $p(x)$ because of sharp boundaries we can use a **smooth kernel**, e.g. a Gaussian

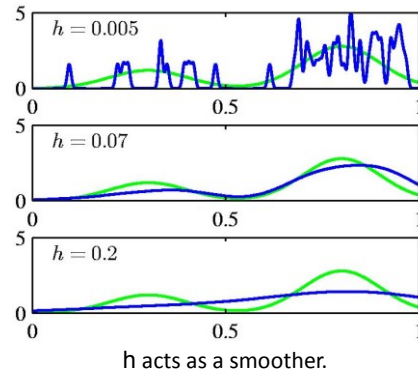
$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{D/2}} \exp\left[-\frac{\|\mathbf{x} - \mathbf{x}_n\|^2}{2h^2}\right]$$

- Any kernel such that

$$k(\mathbf{u}) \geq 0$$

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

- will work.



Nonparametric Methods: kNN estimation

Solution 2: Estimate the probability for \mathbf{x} based on a fixed count \mathbf{K} for a variable volume \mathbf{V} built around \mathbf{x}

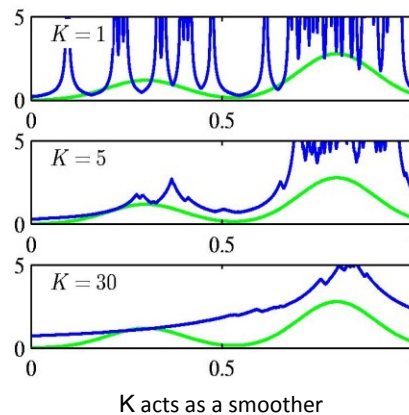
fix \mathbf{K} , estimate \mathbf{V} from the data

Nearest Neighbour Density Estimation:

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

Then

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



Nonparametric vs Parametric Methods

Nonparametric models:

- More flexibility – no density model is needed
- But require storing the entire dataset
- and the computation is performed with all data examples.

Parametric models:

- Once fitted, only parameters need to be stored
 - They are much more efficient in terms of computation
 - But the model needs to be picked in advance
-