

# CS 3750 Machine Learning

## Lecture 3

### Density estimation

Milos Hauskrecht

[milos@cs.pitt.edu](mailto:milos@cs.pitt.edu)

Sennott Square, x4-8845

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

---

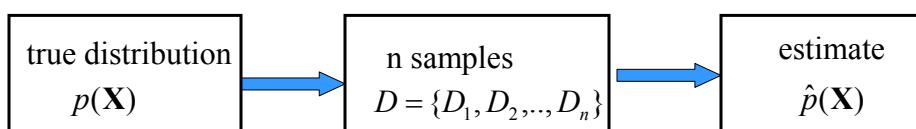
CS 3750 Advanced Machine Learning

### Density estimation

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

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

**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

## Parameter learning

What is the best set of parameters?

- Maximum likelihood (ML) estimates

$$\text{maximize } p(D | \Theta, \xi)$$

$\xi$  - represents prior (background) knowledge

- Maximum a posteriori probability (MAP) estimate

$$\text{maximize } p(\Theta | D, \xi)$$

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

$$P(\mathbf{x} | D, \xi) = \int_{\theta} P(\mathbf{x} | \theta, \xi) p(\theta | D, \xi) d\theta$$

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

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

---

CS 3750 Advanced Machine Learning

## Parameter learning of basic distributions

### Covered :

- Bernoulli
- Binomial
- Multinomial

### Exponential family of distributions

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

- Bernoulli - Beta
- 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:

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

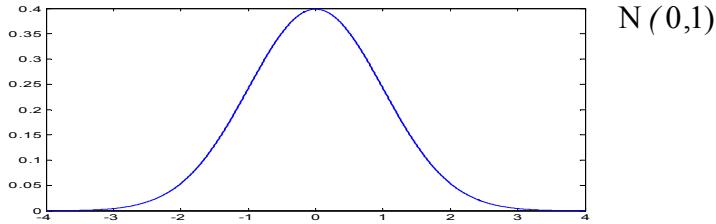
CS 3750 Advanced Machine Learning

## Gaussian (normal) distribution

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



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

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

- **Unbiased estimate:**

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

CS 3750 Advanced Machine Learning

## Multivariate normal distribution

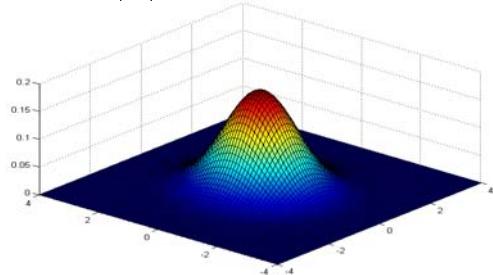
- **Multivariate normal:**  $\mathbf{x} \sim N(\boldsymbol{\mu}, \Sigma)$

- **Parameters:**  $\boldsymbol{\mu}$ - mean  
 $\Sigma$ - covariance matrix

- **Density function:**

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

- **Example:**



CS 3750 Advanced Machine Learning

## Parameter estimates

- **Loglikelihood**

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

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

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

– Covariance estimate is biased

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

- **Unbiased estimate:**

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

CS 3750 Advanced Machine Learning

## Posterior of a multivariate normal

- Assume a prior on the mean  $\mu$  that is normally distributed:

$$p(\mu) \approx N(\mu_p, \Sigma_p)$$

- Then the posterior of  $\mu$  is normally distributed

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

---

CS 3750 Advanced Machine Learning

## Posterior of a multivariate normal

- Then the posterior of  $\mu$  is normally distributed

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

$$\Sigma_n^{-1} = n\Sigma^{-1} + \Sigma_p^{-1}$$

$$\mu_n = \Sigma_p \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \right) + \frac{1}{n} \Sigma \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \mu_p$$

$$\Sigma_n = \Sigma_p \left( \Sigma_p + \frac{1}{n} \Sigma \right)^{-1} \frac{1}{n} \Sigma$$

---

CS 3750 Advanced Machine Learning

## Recursive Bayesian parameter estimation.

- **Recursive Bayesian approach**

- Estimates of the posterior can be sometimes 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

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

- **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 3750 Advanced 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  
$$Z(\boldsymbol{\eta}) = \int h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T t(\mathbf{x}) \} d\mathbf{x}$$
- 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 3750 Advanced 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

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

- Parameters

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

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

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

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

---

CS 3750 Advanced 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

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

- Parameters

$$\boldsymbol{\eta} = \log \frac{\pi}{1 - \pi} \quad (\text{note } \pi = \frac{1}{1 + e^{-\eta}})$$

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

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

$$h(\mathbf{x}) = 1$$

---

CS 3750 Advanced 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**

$$\boldsymbol{\eta} = ? \quad t(\mathbf{x}) = ?$$

$$Z(\boldsymbol{\eta}) = ? \quad h(\mathbf{x}) = ?$$

---

CS 3750 Advanced 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**

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

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

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

---

CS 3750 Advanced Machine Learning

## Exponential family

- For iid samples, the likelihood of data is

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

- Important:

- the dimensionality of the sufficient statistic remains the same with the number of samples

---

CS 3750 Advanced Machine Learning

## Exponential family

- 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 3750 Advanced Machine Learning

## Exponential family

- Rewriting the gradient:

---

CS 3750 Advanced Machine Learning

## Exponential family

- Rewriting the gradient:

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

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

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

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

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

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

---

CS 3750 Advanced Machine Learning

## Moments of the distribution

- **For the exponential family**

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

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

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

- **Derivatives:**

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

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

---

CS 3750 Advanced Machine Learning