

STAT 391

4/27/23

lecture 10

- Gradient ascent
- Non-parametric density estimation
(kernel d. e.)

LV non-parametric

Lecture Notes IV – Continuous distributions. Parametric density estimation.

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

April, 2023

CDF and PDF sampling



Examples of continuous distributions



ML estimation for continuous distributions



ML estimation by gradient ascent



Reading: Ch.5, 6

Examples of continuous distributions

$$\mathcal{F}_1 = \{u_{[a,b]}, a < b\} \quad \text{uniform}$$

(5)

$$f(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$

(6)

$$\mathcal{F}_2 = \{N(\cdot; \mu, \sigma^2)\} \quad \text{normal}$$

(7)

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

(8)

$$F(x; a, b) = \frac{1}{1 + e^{-ax-b}}, \quad a > 0 \quad \text{logistic}$$

(9)

$$f(x; a, b) = \frac{ae^{-ax-b}}{(1 + e^{-ax-b})^2}$$

(10)

ML estimation by gradient ascent

Maximize $l(a, b)$ by gradient ascent

Algorithm

Initialize a^0, b^0 $a > 0$

Until convergence

$t = 1, 2, \dots$

$$a^t \leftarrow a^{t-1} + \eta \frac{\partial l}{\partial a}(a^t, b^t)$$

$$b^t \leftarrow b^{t-1} + \eta \frac{\partial l}{\partial b}(a^t, b^t)$$

Output $\hat{a}^t, \hat{b}^t = a^M, b^M$

Stopping G.A.

a) always have t_{\max} limit

b) relative increase in l

$$10^{-10} < \text{tolerance} \sim \frac{1}{\sqrt{n}} \left(\frac{1}{100} \right)$$

$$\frac{l(a^t, b^t) - l(a^{t-1}, b^{t-1})}{l(a^{t-1}, b^{t-1})} \leq \text{tolerance}$$

ML estimation by gradient ascent

$$l(a, b) = n \ln a - a \sum_i x_i - nb - 2 \sum_{i=1}^n \ln(1 + e^{-ax_i - b})$$

Normalize $\ell(a, b)$!! ❤

$$\ell^{norm}(a, b) \leftarrow \frac{\ell(a, b)}{n}$$

$$\frac{\partial l}{\partial a} = \frac{n}{a} - \sum_{i=1}^n x_i \frac{1 - e^{-ax_i - b}}{1 + e^{-ax_i - b}} = 0$$

$$\frac{\partial l}{\partial b} = - \sum_{i=1}^n \frac{1 - e^{-ax_i - b}}{1 + e^{-ax_i - b}} = 0$$

3) $\sqrt{\left(\frac{\partial \ell}{\partial a}\right)^2 + \left(\frac{\partial \ell}{\partial b}\right)^2} < tol$

length of $\begin{bmatrix} \frac{\partial \ell}{\partial a} \\ \frac{\partial \ell}{\partial b} \end{bmatrix}$

ML estimation by gradient ascent

$$l(a, b) = n \ln a - a \sum_i x_i - nb - 2 \sum_{i=1}^n \ln(1 + e^{-ax_i - b})$$

$$\frac{\partial l}{\partial a} = \frac{n}{a} - \sum_{i=1}^n x_i \frac{1 - e^{-ax_i - b}}{1 + e^{-ax_i - b}} = 0$$

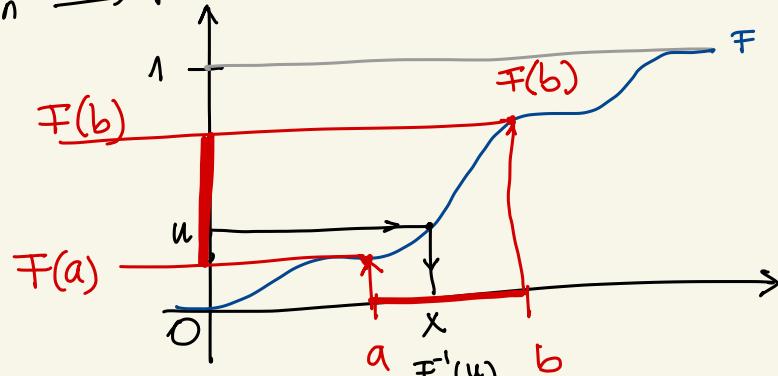
$$\frac{\partial l}{\partial b} = - \sum_{i=1}^n \frac{1 - e^{-ax_i - b}}{1 + e^{-ax_i - b}} = 0$$

Sampling from a continuous distribution

Given F CDF of P on $(-\infty, \infty)$

Wanted sample $x \sim P \Leftrightarrow$ if many samples taken $\rightarrow P$

1. sample $u \sim U[0, 1]$
 $u = \text{rand();}$
2. output $x = F^{-1}(u)$



Why?

$$\Pr[a, b] = F(b) - F(a)$$

↑
Alg

$$\Pr[x \in [a, b]]$$

$$\Pr[x \in [a, b]] = \Pr[u \in [F(a), F(b)]] = F(b) - F(a)$$

$$x = F^{-1}(u) \Leftrightarrow F(x) = u$$

definition of uniform $[0, 1]$

Lecture Notes V – Non-parametric density estimation

S discrete ✓
ML

S continuous
 $(-\infty, \infty)$

models ↗ parametric ✓ ML

$\mathcal{F} = \{f_\theta\}$ ↗ non-parametric ← NO ML !!

April, 2023

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

Kernels

Kernel density estimators

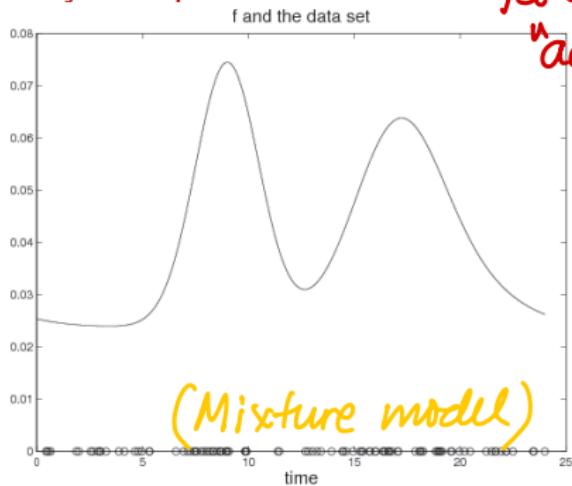
Choosing h by Cross-Validation and the Bias-Variance trade-off

Reading: Ch. 7

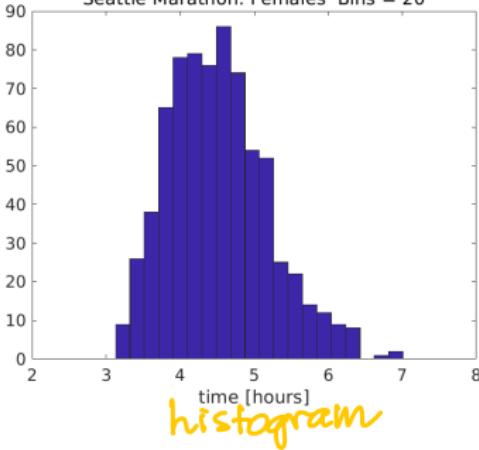
Why non-parametric?

— fit any shape

"adapt" to shape of data



Seattle Marathon: Females Bins = 20



Parametric families

Normal



Logistic



Exponential



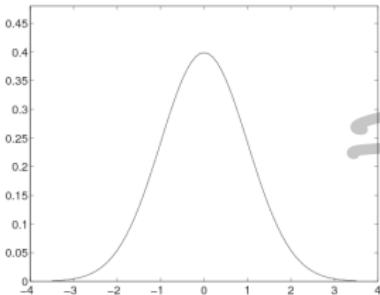
[Cauchy]



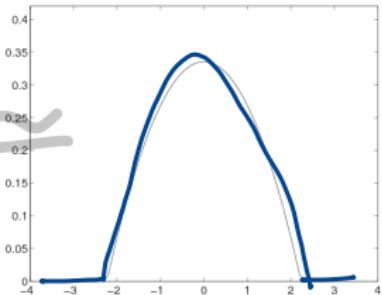
Kernel density estimation ↙

Kernels

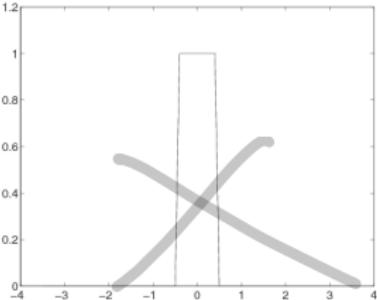
(Shape of kernel not important)



Gaussian
kernel
 $N(0,1)$



Epanechnikov



Square
unif $[-0.5, 0.5]$

kernel = ^{simple} density $\rightarrow k(z) \geq 0$
 $\int_{-\infty}^{\infty} k(z) dz = 1$

$$k(z)$$

$$k: (-\infty, \infty) \rightarrow [0, \infty)$$

$$1. \quad k(-z) = k(z)$$

$$2. \quad k(z) \rightarrow 0 \text{ for } z > 0$$

$$[3.] \quad \int_{-\infty}^{\infty} z^2 k(z) dz = 1 \text{ standardization}$$

Kernel density estimators

$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n k(\mathbf{x} - \mathbf{x}_i)$$

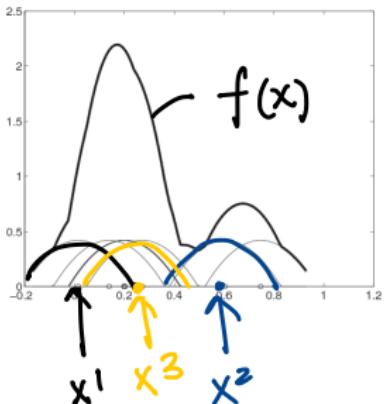
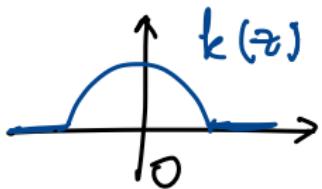
data

normalization

$$k(x - x_i)$$

$\xrightarrow{\text{shift to } x_i}$

\tilde{z}

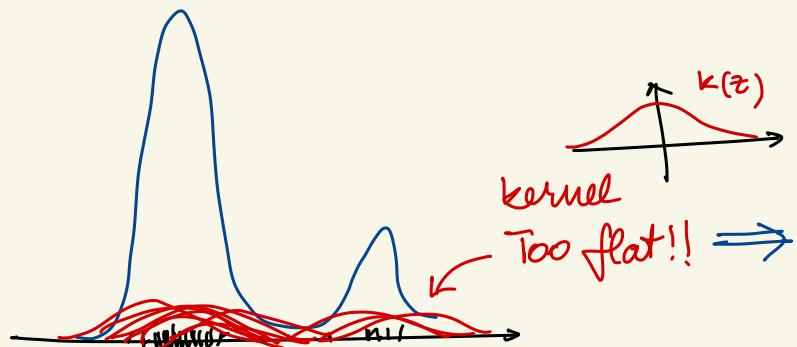


$$\int_{-\infty}^{\infty} f(x) dx = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} k(x - x_i) dx$$

$\underbrace{\hspace{100px}}_1$

$$= \frac{1}{n} \cdot n = 1$$

\Downarrow
f is a density



(bandwidth)
 $h = \text{kernel width}$

$$k_h(z) = \frac{1}{h} k\left(\frac{z}{h}\right)$$

↑
unit on
x axis
(z)

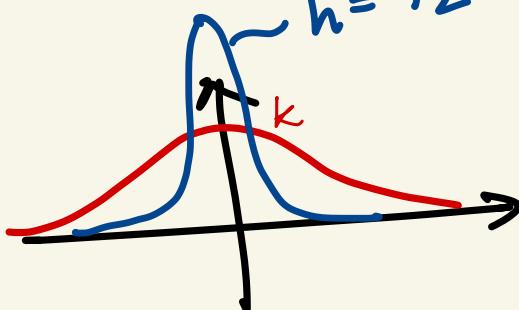
$$f'(x) = \frac{1}{n} \sum_i k'(x - x_i) \Rightarrow |f'| \leq \frac{1}{n} \sum_{i=1}^n a = a$$

$$|k'(z)| < a$$

$$h = \frac{1}{2}$$

Exercise:

$$\int_{-\infty}^{\infty} k_h(z) dz = 1$$



Kernel density estimators

$$f_h(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right)$$

↑ width

data

smoothness
param

