

STAT 391

4/18/23

Lecture 8

Quiz 1
on 4/20
at 12:30

Smoothing example
Continuous S +
parametric

Sol 1, 2
L1V posted

Lecture Notes III: Discrete probability in practice – Small Probabilities

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

April, 2023

Definitions and setup

We will look at estimating categorical distributions from samples, when the number of outcomes m is large.

- ▶ Let $S = \{1, \dots, m\}$ be the sample space, and $P = (\theta_1, \dots, \theta_m)$ a distribution over S .
- ▶ We draw n independent samples from P , obtaining the **data set** \mathcal{D}
- ▶ Define the **counts** $\{n_j = \#\# j \text{ appears in } \mathcal{D}, i = 1, \dots, n\}$. The counts are also called **sufficient statistics** or **histogram**.
- ▶ Define the **fingerprint** (or **histogram of histogram**) of \mathcal{D} as the counts of the counts, i.e $\{r_k = \#\text{counts } n_j = k, \text{ for } k = 0, 1, 2, \dots\}$

Example $m = 26$ alphabet letters

Data

the red fox is quick
 $n = 16$ letters

ho ho who s on first
 $n = 15$ letters

Counts n_i

$n_j = 0 : a, b, g, j, l, m, n,$
 $n_j = 1 : v, w, y, z$
 $n_j = 2 : c, d, f, h, k, o, q, r, s, t, u, x$
 $n_j = 2 : e, i$

$n_j = 0 : a, b, c, \dots, x, z$
 $n_j = 1 : f, i, n, r, t, w$
 $n_j = 2 : s$
 $n_j = 3 : h$
 $n_j = 4 : o$

Fingerprint r_k

$r_0 = 12 = |\{a, b, g, \dots, y, z\}|$
 $r_1 = 12 = |\{c, d, f, h, \dots, u, x\}|$
 $r_2 = 2 = |\{e, i\}|$
 $r_3 = \dots r_n = 0$

$r_0 = 26 - 6 - 1 - 1 - 1 = 17$
 $r_1 = 6 = |\{f, i, n, r, t, w\}|$
 $r_2 = 1 = |\{s\}|$
 $r_3 = 1 = |\{h\}|$
 $r_4 = 1 = |\{o\}|$

- ▶ It is easy to verify that $n_j \in 0 : n$, hence $r_{0:n}$ may be non-zero (but $r_{n+1,n+2,\dots} = 0$), and that

$$m = r_0 + r_1 + \dots + r_n \quad n = 0 \times r_0 + 1 \times r_1 + \dots + k \times r_k + \dots \quad (1)$$

Smoothing on an example

$$(a)_+ = \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

 $n = \text{sample size}$ $m = \text{alphabet size} = |\Sigma|$

- the counts $\{n_j = \#j \text{ appears in } D, i = 1, \dots, n\}$ (or **sufficient statistics** or **histogram**)
- fingerprint (or **histogram of histogram**) of D as the counts of the counts $\{r_k = \#\text{counts } n_j = k, \text{ for } k = 0, 1, 2, \dots\}$, and $R_k = \{j, n_j = k, \}$

$$m = r + r_0$$

\uparrow
obs

\uparrow
lunohs

Example $m = 26$ alphabet letters

Data

Counts n_i

$$n_i = 0: a, b, g, j, l, m, n,$$

$$p, v, w, y, z$$

$$n_i = 1: c, d, f, h, k, o, q, r, s, t, u, x$$

$$n_i = 2: e, i$$

$$\Theta^{ML} = 0$$

a,b

Fingerprint r_k

$$r_0 = 12 = |\{a, b, g, \dots, y, z\}|$$

$$r_1 = 12 = |\{c, d, f, h, \dots, u, x\}|$$

$$r_2 = 2 = |\{e, i\}|$$

$$r_3 = \dots r_n = 0$$

the red fox is quick

 $n = 16$ letters

$$n_j = 0$$

$$n_j > 0$$

$$r = r_1 + r_2 + \dots$$

$$(\text{Lap}) \quad n_j \leftarrow n_j + 1$$

$$\text{NE} \quad \Theta_j^{\text{NE}} = \frac{r/m}{n}$$

WB

$$P[\text{NEW}] = \frac{r}{n}$$

$$\Theta_j^{\text{WB}} = \frac{1}{r_0} \frac{r}{n}$$

$$\text{GT} \quad P_1 = \frac{r_1}{n}$$

$$\Theta_j^{\text{GT}} = \frac{1}{r_0} \frac{r_1}{n}$$

$$\Theta_j^{\text{NE}} = \frac{(n_j - 1) + r/m}{n}$$

NE
 n_j

CORRECTION!!

~~$$\Theta_j^{\text{WB}} = \frac{n_j}{n} \left(1 - \frac{r}{n}\right) \frac{1}{1 + \frac{r}{n}}$$~~

$$\Theta_j^{\text{GT}} = \frac{n_j}{n} \left(1 - \frac{r_1}{n}\right)$$

Smoothing on an example

- ▶ the counts $\{n_j = \#j \text{ appears in } \mathcal{D}, i = 1, \dots, n\}$ (or **sufficient statistics** or **histogram**)
- ▶ fingerprint (or **histogram of histogram**) of \mathcal{D} as the counts of the counts $\{r_k = \#\text{counts } n_j = k, \text{ for } k = 0, 1, 2, \dots\}$, and $R_k = \{j, n_j = k, \}$

Example $m = 26$ alphabet letters

Data

the red fox is quick
 $n = 16$ letters

Counts n_i

$n_j = 0: a, b, g, j, l, m, n,$

p, v, w, y, z

$n_j = 1: c, d, f, h, k, o, q, r, s, t, u, x$

$n_j = 2: e, i$

Fingerprint r_k

$r_0 = 12 = |\{a, b, g, \dots, y, z\}|$

$r_1 = 12 = |\{c, d, f, h, \dots, u, x\}|$

$r_2 = 2 = |\{e, i\}|$

$r_3 = \dots r_n = 0$

$$n_j = 0$$

$$\Theta_j^{NE} = \frac{r/m}{n} = \frac{7}{13 \cdot 16}$$

$$\Downarrow r = m - r_0 = 26 - 12 = 14$$

$$r/m = \frac{14}{26} = \frac{7}{13}$$

$$\Theta_j^{WB} = \frac{1}{r_0} \frac{r}{n} = \frac{7}{8} \frac{1}{15} \frac{1}{12} > \Theta_j^{NE}$$

$$\frac{r}{n} = \frac{14}{16} = \frac{7}{8}$$

$$\Theta_j^G = \frac{1}{r_0} \frac{r}{n} = \frac{3}{4} \cdot \frac{1}{12} = \frac{1}{16} < \Theta_j^{WB}$$

$$\text{always } \frac{r}{n} = \frac{12}{16} = \frac{3}{4}$$

$$\frac{r/n}{1+r/n} = \frac{r}{n+r} = \frac{14}{30} = \frac{7}{15}$$

Smoothing on an example

- the counts $\{n_j = \#j \text{ appears in } \mathcal{D}, i = 1, \dots, n\}$ (or **sufficient statistics** or **histogram**)
- fingerprint (or **histogram of histogram**) of \mathcal{D} as the counts of the counts $\{r_k = \#\text{counts } n_j = k, \text{ for } k = 0, 1, 2, \dots\}$, and $R_k = \{j, n_j = k, \}$

Example $m = 26$ alphabet letters

Data

the red fox is quick
 $n = 16$ letters

Counts n_i

$$n_j = 0: a, b, g, j, l, m, n,$$

$$p, v, w, y, z$$

$$n_j = 1: c, d, f, h, k, o, q, r, s, t, u, x$$

$$n_j = 2: e, i$$

Fingerprint r_k

$$r_0 = 12 = |\{a, b, g, \dots, y, z\}|$$

$$r_1 = 12 = |\{c, d, f, h, \dots, u, x\}|$$

$$r_2 = 2 = |\{e, i\}|$$

$$r_3 = \dots r_n = 0$$

$$r = m - r_0 = 26 - 12 = 14$$

$$\theta_j^{\text{NE}} = \frac{(n_j - 1) + r/m}{n} \Rightarrow \theta_{c,d,\dots}^{\text{NE}} = \frac{7}{13/16}$$

$$\theta_{e,i}^{\text{NE}} = \frac{1 + 7/13}{16}$$

$$\theta_j^{\text{WB}} = \frac{n_j}{n} \left(1 - \frac{r}{n}\right)$$

$$\theta_{c,d,\dots}^{\text{WB}} = \frac{1}{16} \cdot \frac{1}{8} \cancel{\frac{8}{15}}$$

$$\theta_j^{\text{GT}} = \frac{n_j}{n} \left(1 - \frac{r_1}{n}\right)$$

$$\theta_{e,i}^{\text{WB}} = \frac{2}{16} \cdot \frac{1}{8} \cancel{\frac{8}{15}}$$

$$\frac{r}{n} = \frac{14}{16} - \frac{7}{8}$$

$$r_0 = 12$$

$$\frac{r_1}{n} = \frac{12}{16} = \frac{3}{4}$$

$$\frac{1}{1 + \frac{r}{n}} = \frac{n}{n+r} = \frac{16}{30} = \frac{8}{15}$$

$$1 - \frac{r}{n} = \frac{1}{8}$$

(Estimating)

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

CDF and PDF refresher

$S = (-\infty, \infty)$ or uncountable

Cumulative distribution function (CDF)

$$F(x) = P[X \leq x]$$

$$= \Pr[-\infty, x] \text{ as } (1)$$

→ 1. $F \geq 0$ positivity.

→ 2. $\lim_{x \rightarrow -\infty} F = 0 = \Pr[-\infty]$

→ 3. $\lim_{x \rightarrow \infty} F = 1 = \Pr(-\infty, \infty) = 1$

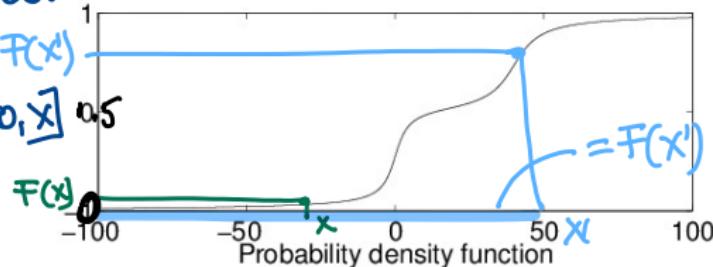
→ 4. F is an increasing function

Probability density [function]
(PDF)

$$f = \frac{dF}{dx} \quad (2)$$

e.g. $S = (a, b)$

Cumulative distribution function



$$P(a, b) = \underline{P[a, b]} = \underline{F(b) - F(a)} = \int_a^b f(x) dx \Rightarrow F(b) \geq F(a) \quad (3)$$

normalization condition

$$\int_{-\infty}^{\infty} f(x) dx = 1 \quad (4)$$

CDF and PDF refresher

Cumulative distribution function (CDF)

$$F(x) = P[X \leq x] \quad (1)$$

1. $F \geq 0$ positivity.
2. $\lim_{x \rightarrow -\infty} F = 0$
3. $\lim_{x \rightarrow \infty} F = 1$
4. F is an increasing function

Probability density [function] (PDF)

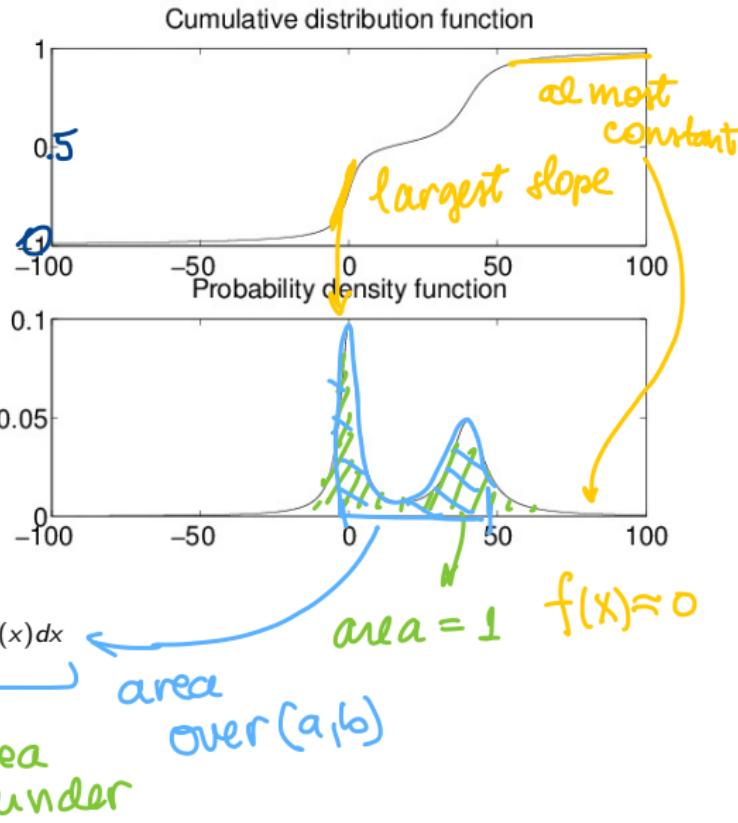
$$f = \frac{dF}{dx} \quad (2)$$

$$P(a, b) = P[a, b] = F(b) - F(a) = \int_a^b f(x) dx \quad (3)$$

normalization condition

$$2. \int_{-\infty}^{\infty} f(x) dx = 1 \quad (4)$$

1. $f(x) \geq 0$ for all x $f(x)$



Notation

$$P[\mathcal{E}]$$

↑ event

$$[a, \infty), (-\infty, b), \dots$$

in $(-\infty, \infty)$, events are $[a, b], (a, b), [a, b], \dots$
or unions of intervals

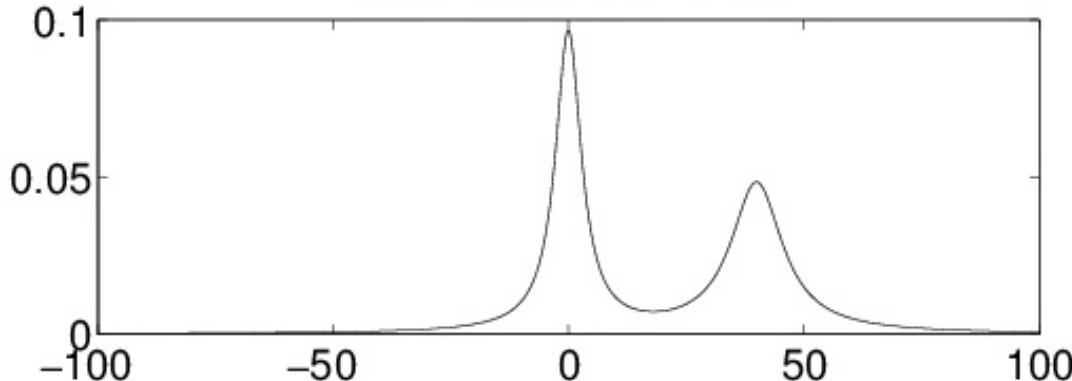
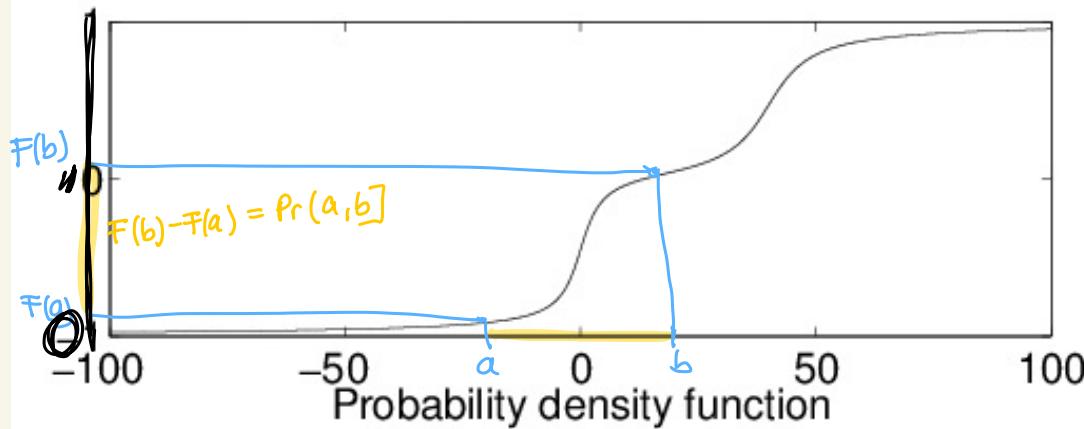
$$P[a, b] = P[a, b]$$

simplification

Assume 1) $F(x)$ continuous $\Rightarrow P[a] = P[a, a] = F(a) - F(a) = 0$

2) $F(x)$ differentiable $\Rightarrow f(x)$ exists

Cumulative distribution function



Examples of continuous distributions

family of distributions

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

$$f(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

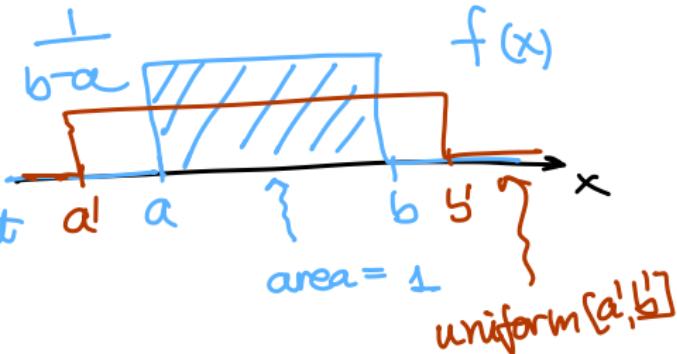
$$\mathcal{F}_2 = \{N(\cdot; \mu, \sigma^2)\} \quad \text{normal} \quad (6)$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (7)$$

$$F(x; a, b) = \frac{1}{1 + e^{-ax-b}}, \quad a > 0 \quad \text{logistic} \quad (8)$$

$$f(x; a, b) = \frac{ae^{-ax-b}}{(1 + e^{-ax-b})^2} \quad (9)$$

\mathcal{F}_1 has parameters a, b , $a < b$



~~uniform $(-\infty, \infty)$~~

$$\int_{-\infty}^{\infty} c dx = \infty$$

for any $c > 0$

Examples of continuous distributions

Normal: $\mu = \text{mean} \in (-\infty, \infty)$
 $\sigma^2 = \text{variance} > 0$ } parameters for \mathcal{F}_2

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

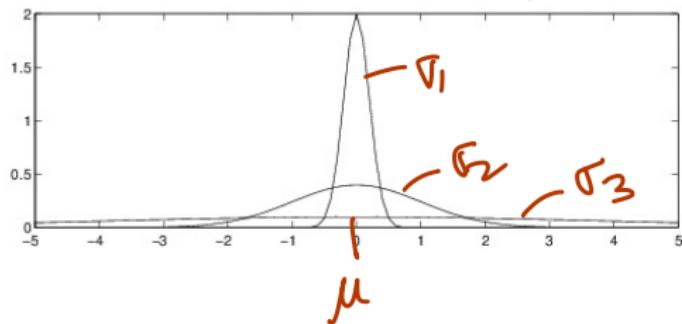
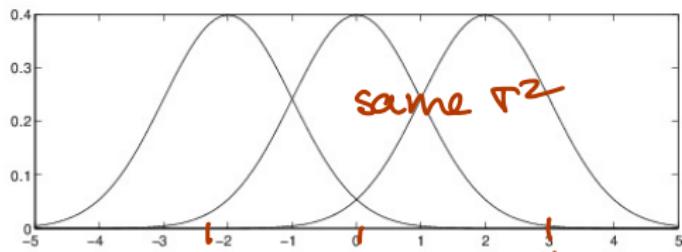
$$f(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

$$\mathcal{F}_2 = \{N(\cdot; \mu, \sigma^2)\} \quad \text{normal} \quad (7)$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (8)$$

$$F(x; a, b) = \frac{1}{1 + e^{-ax-b}}, a > 0 \quad \text{logistic} \quad (9)$$

$$f(x; a, b) = \frac{ae^{-ax-b}}{(1 + e^{-ax-b})^2} \quad (10)$$



$$\sigma_3 > \sigma_2 > \sigma_1$$

Examples of continuous distributions

a, b parameters for
 $a > 0$ \mathcal{F}_3

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

(5)

$$f(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

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

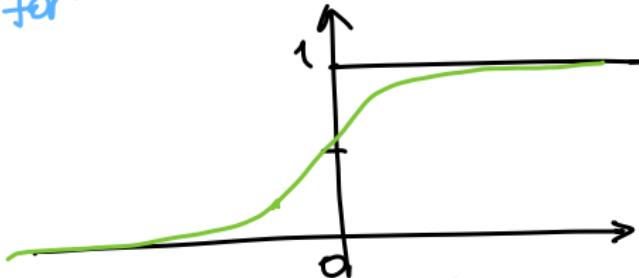
(7)

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (8)$$

$$F(x; a, b) = \frac{1}{1 + e^{-ax-b}}, a > 0 \text{ logistic} \quad (9)$$

$$f(x; a, b) = \frac{ae^{-ax-b}}{(1 + e^{-ax-b})^2} \quad (10)$$

Logistic



$$a=1 \Rightarrow F(x) = \frac{1}{1+e^{-x}}$$

$$x=0 \Rightarrow F(x) = \frac{1}{2}$$

$$x \rightarrow \infty \Rightarrow F(x) \rightarrow 1$$

$$x \rightarrow -\infty \Rightarrow F(x) \rightarrow 0$$

$$F(x) = 1 - F(-x)$$

Exercise ↗

Examples of continuous distributions

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

(5)

$$f(x; a, b) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases} \quad (6)$$

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

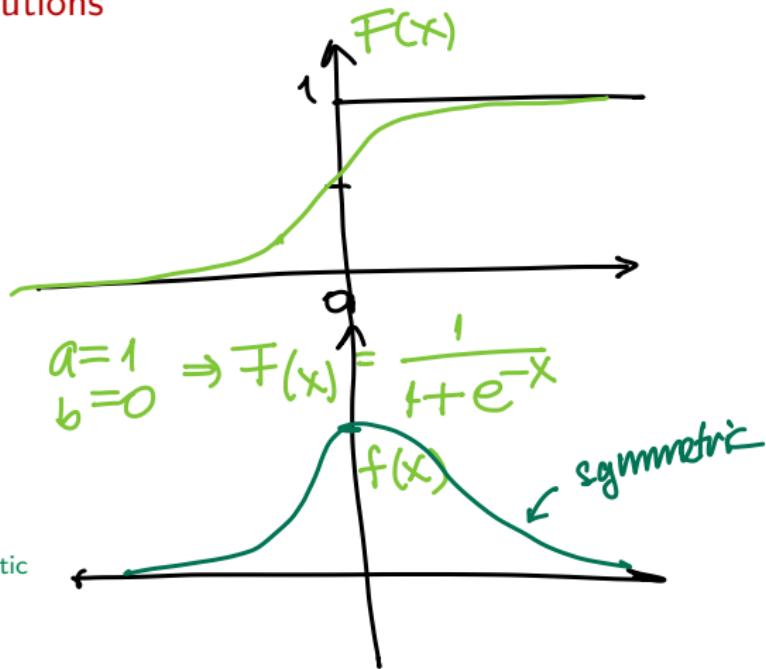
(7)

$$f(x; \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (8)$$

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

(9)

$$f(x; a, b) = \frac{ae^{-ax-b}}{(1 + e^{-ax-b})^2} \quad (10)$$

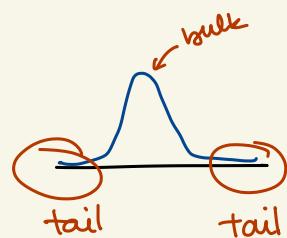


$$f(x) = f(-x) \quad \text{for } a=1, b=0$$

$$\arg\max_x f(x) \Leftrightarrow ax+b=0 \Leftrightarrow x=-\frac{b}{a}$$

About "tails" regions

$$N(0, 1) \rightarrow f_1(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}$$



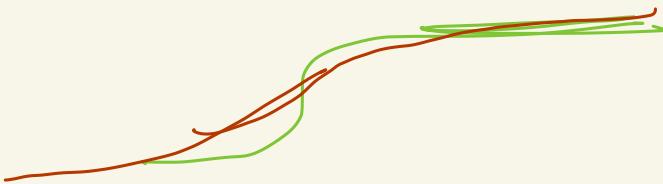
$$\text{logistic } (a=1, b=0) \rightarrow f_2(x) = \frac{1}{(1+e^{-x})^2}$$

$$(x \rightarrow \infty) \quad \underline{\text{Large}} \quad \approx 1$$

$$\frac{f_1}{f_2} \approx \frac{1}{\sqrt{2\pi}} \frac{e^{-\frac{x^2}{2}}}{e^{-x}} = \frac{1}{\sqrt{2\pi}} e^{x - \frac{x^2}{2}} \rightarrow 0 !!$$

$$x \underline{\text{large}} \quad x - \frac{x^2}{2} \rightarrow -\infty$$

$$f_1(x) < f_2(x) \quad \text{for } |x| \text{ large}$$



ML estimation for continuous distributions = density estimation

$$S = (-\infty, \infty) \quad \text{or} \quad S \subset (-\infty, \infty)$$

Sample space

Data $\mathcal{D} = \{x^1, x^2, \dots, x^n\} \subset S$ iid from true unknown f

Model family $\mathcal{F} = \{f(x|\theta), \theta = \text{parameters}\}$

Chosen \uparrow Normal, Uniform, ...

Problem estimate θ from \mathcal{D}

Solution : Max likelihood $\rightarrow \theta^{ML} \rightarrow f(x|\theta^{ML})$
estimated \uparrow
density

ML estimation for continuous distributions

Likelihood of θ

$$L(\theta) = \prod_{i=1}^n f(x^i | \theta)$$

$P(x^i | \theta)$ for S discrete

$F(x)$

log-likelihood

$$\ell(\theta) = \sum_{i=1}^n \ln f(x^i | \theta)$$

Max Likelihood Principle

$$\text{choose } \theta^{ML} = \operatorname{argmax} \ell(\theta)$$

-

x^i observed

$$P[x^i] = 0$$

$$P[x^i - \Delta, x^i + \Delta] \cong 2\Delta \cdot f(x^i)$$

$$L(\theta) = \prod_i P(x^i) \cong \prod_i 2\Delta f(x^i) = (2\Delta)^n \prod_i f(x^i | \theta)$$

↳ independent
of θ

Deriving $L(\theta)$

by analogy with discrete case