

# Lecture 8

Small probs  
Continuous S

Happy Lunar New  
Year!

Q1 Thu 2/6  
LIV Cont distributions  
HW2 due  
Sol 2 Tue 2/4  
HW3 TR posted

# Lecture Notes III: Discrete probability in practice – Small Probabilities

Marina Meilă  
mmp@stat.washington.edu

Department of Statistics  
University of Washington

January, 2025

The problem with estimating small probabilities ✓

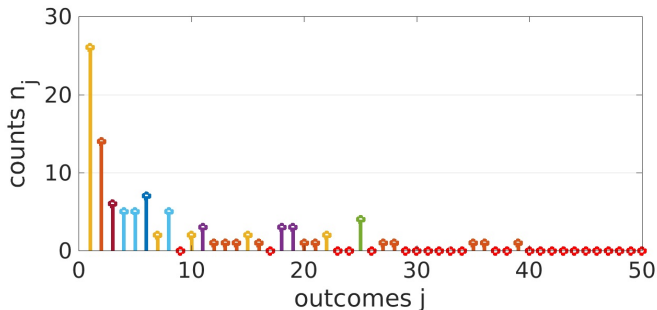
Definitions and setup ✓

Additive methods (Laplace, Dirichlet, Bayesian, ELE) ✓

Discounting (Ney-Essen) ✓

Multiplicative smoothing: Estimating the next outcome (Witten-Bell, Good-Turing) ↙

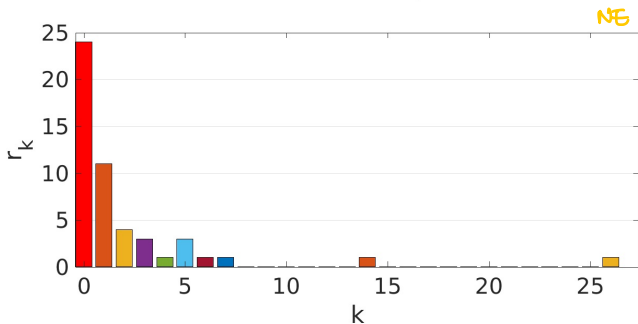
~~Back-off or shrinkage – mixing with simpler models~~



$R_k = \{j \text{ that appear } k \text{ times}\}$

$$n_j = k$$

$$\Theta_j^{ML} = \frac{k}{n} \rightarrow \tilde{\theta}_j$$



NE

$$R_0 \cup R_1: n_j^i = 0$$

$$R_k \quad k > 1$$

WB  $R_0$  special  
 $R_k \quad k > 1$

GT  $R_k$  separate  
for all  $k$

## 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**  $n_j$

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

$p, v, w, y, z$

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

$n_j = 2 : e, i$

$n_j = 0 : a, b, c, \dots, x, z$

$n_j = 1 : f, i, n, r, t, w$

$n_j = 2 : s$

$n_j = 3 : h$

$n_j = 4 : o$

**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

$$\text{Extra } \tilde{\theta}_{R_1}^{WB} < \tilde{\theta}_{R_0}^{WB} ?$$

- ▶ 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  $n_j$

$n_j = 0$ : 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  $r_k$

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

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

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

$r_3 = \dots r_n = 0$

$$r_2 = m - r_0 = 14$$

$$m = 12 + 14 = 26$$

Witten-Bell

$y_i = 1$  if  $x^i$  "new"  
 0 otherwise

$$p_0 = \Pr[y_i = 1] = \frac{\# \text{new}}{n} = \frac{r_2}{n}$$

$$j \in R_0 \Rightarrow \Pr[X^{n+1} = j] = \frac{p_0}{r_0}$$

$$n_j > 0 : \theta_j^{ML} = \frac{n_j}{n} \rightarrow \sum_j \theta_j^{ML} = 1$$

$$n_j = 0 : \theta_j = \frac{p_0}{r_0} \rightarrow \sum_{j \in R_0} \theta_j = p_0$$

$\rightarrow 1 + p_0 = Z$  normalization

• no sense if  $r_2 \lesssim m$

• some sense  $r_2 \ll m$

$n$  large  $\gg m \gg r_2$   
 want  $\Pr[X^{n+1} \text{ is new}] \approx \frac{r_2}{n}$

# 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**  $n_j$

$n_j=0$ : 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**  $r_k$

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

$$\begin{aligned}
 n_j > 0: \theta_j^{ML} &= \frac{n_j}{n} \rightarrow \sum_j \theta_j^{ML} = 1 \rightarrow \tilde{\theta}_j^{WB} = \frac{\theta_j^{ML}}{1+p_0} \leftarrow \text{rescale } \frac{1}{1+p_0} \text{ by } \frac{1}{1+p_0} \\
 n_j = 0: \theta_j &= \frac{p_0}{r_0} \rightarrow \sum_{j \in R_0} \theta_j = p_0 \rightarrow \tilde{\theta}_j^{WB} = \frac{p_0}{r_0} \cdot \frac{1}{1+p_0} \xrightarrow{\approx} \frac{p_0}{1+p_0}
 \end{aligned}$$

TO CHECK:  $\sum_{j=1}^m \tilde{\theta}_j^{WB} = 1$

- when not to use WB

$$\frac{n}{m} \approx 1 !!$$

# Witten-Bell discounting – probability of a new value

## ► Idea:

- Look at the sequence  $(x_1, \dots, x_n)$  as a binary process: either we observe a value of  $X$  that was observed before, or we observe a new one.
- Assume that of  $m$  possible values  $r$  were observed (and  $m - r$  unobserved)
- Then the probability of observing a new value is  $p_0 = \frac{r}{n}$ .
- Hence, set the probability of all unseen values of  $X$  to  $p_0$ . The other probability estimates are renormalized accordingly.

$$\theta_j^{WB} = \begin{cases} \frac{n_j}{n} \frac{1}{1+p_0} = \frac{n_j}{n+r} & n_j > 0 \\ \frac{1}{m-r} \frac{p_0}{1+p_0} = \frac{1}{m-r} \frac{r}{n+r} & n_j = 0 \end{cases} \quad (7)$$

Witten-Bell makes sense only when some  $n_j$  counts are zero. If all  $n_j > 0$  then W-B smoothing has undefined results.

WB smoothing has no parameter to choose (GOOD!)



# Comparison of the methods ← toy ex.

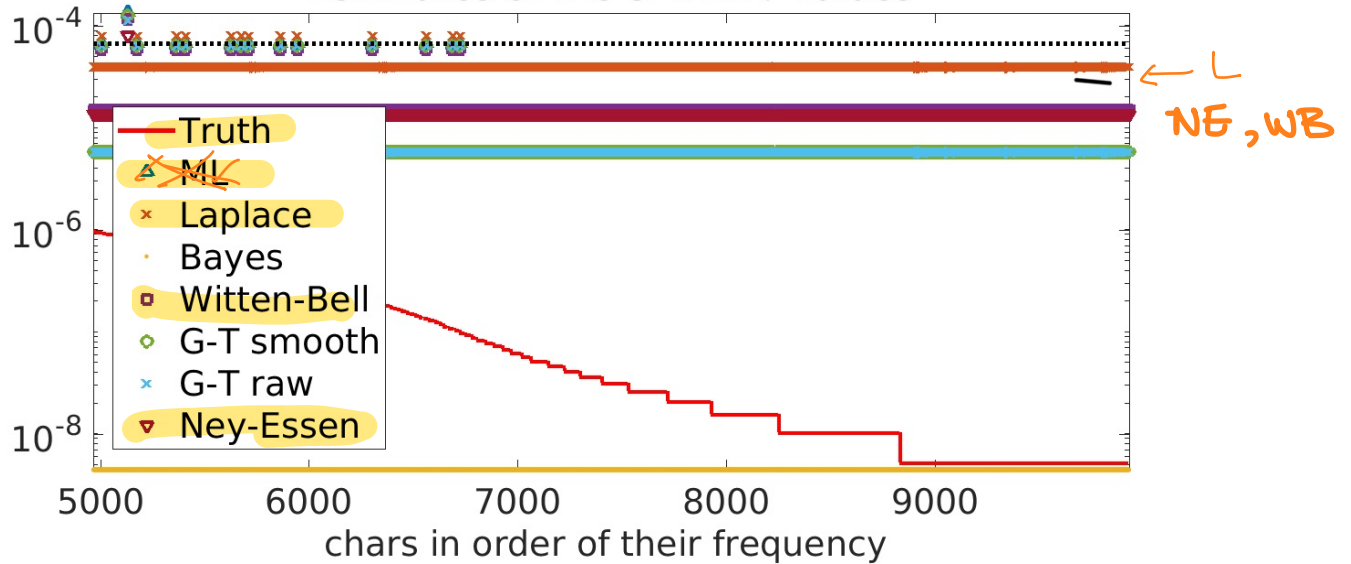
Numerical values to exemplify the results:  $n = 1000$ ,  $m = 1000$ ,  $r = 100$

Count $n_j$	$n_j = 0$	$n_j = 1$	$n_j \gg 1$
$\theta_j^{ML}$	0	$\frac{1}{n} = \frac{1}{1000}$	$\frac{n_j}{n}$
$\theta_j^{Laplace}$	$\frac{1}{n+m} = \frac{1}{2000}$	$\frac{2}{n+m} = \frac{1}{1000}$	$\frac{n_j+1}{n+m} = \frac{n_j+1}{2000} = 0.5 \theta_j^{ML}!!$
$\theta_j^{Bayes}$ , $n^0 = 1$ , $n_j^0 = \frac{1}{m}$	$\frac{1}{m(n+1)} \approx \frac{1}{10^6}$	$\frac{1+1/m}{n+1} \approx \frac{1}{10^3}$	$\frac{n_j+1/m}{n+1} \approx \frac{n_j}{1000}$
$\theta_j^{NE}$ , $r = 1$ = tax	$\frac{r}{mn} = \frac{1}{10^4}$	$\frac{r}{mn} = \frac{1}{10^4}$	$\frac{n_j-1+r/m}{n+1} \approx \frac{n_j}{1000} \leq \theta_j^{ML}$
$\theta_j^{WB}$	$\frac{1}{m-r} \frac{r}{n+r} = \frac{1}{9900}$	$\frac{1}{n+r} = \frac{1}{1100}$	$\frac{n_j}{n+r} = \frac{n_j}{1100}$

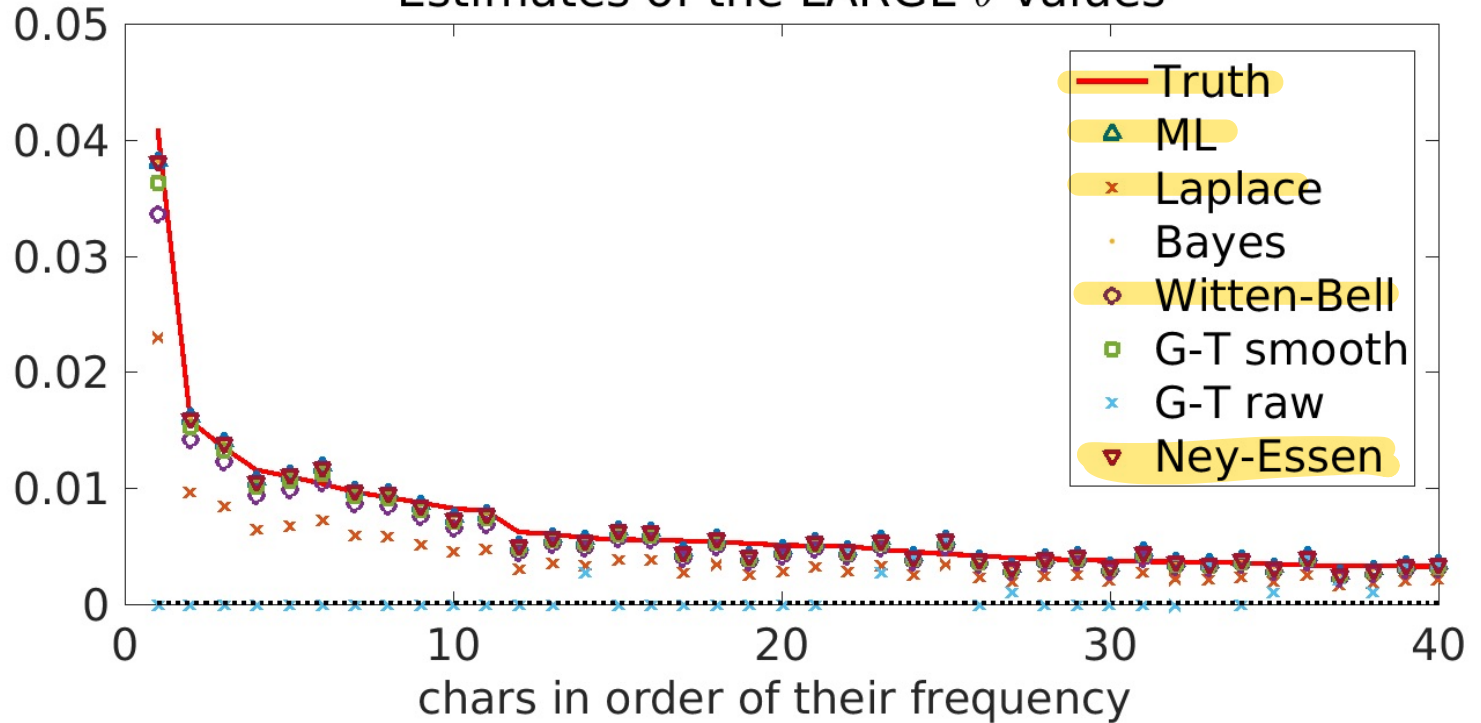
## Remarks

- ▶ Laplace shrinks ML estimates of large probabilities by factor of 2. Too much! (because large  $\theta_j^{ML}$  are close to their true values)
- ▶ Bayes (with uninformative prior) affects large  $\theta_j^{ML}$  much less than small ones. Good
- ▶ Ney-Essen smooths more when  $r$  is larger; any  $n_j$  is affected by less than  $\delta$ .
- ▶ Ney-Essen estimates of  $\theta_j^{NE}$  for counts of 0 and 1 are equal to a fraction of  $\frac{r}{m}$  (this grows with  $n$  as  $r$  grows with  $n$ ).
- ▶ In Witten-Bell, the large  $\theta_j^{ML}$  are shrunk depending on  $r$ , but independently of  $m$ . Proportional, bad
- ▶ ... but, if we overestimate  $m$  grossly, the overestimation will only affect the  $\theta_j^{WB}$  for the 0 counts, but none of the  $\theta_j^{WB}$  for the values observed. (true for NE as well).

Estimates of the SMALL  $\theta$  values



## Estimates of the LARGE $\theta$ values



sample spaces

Lecture Notes IV – Continuous distributions. Parametric density estimation.

Marina Meilă  
mmp@stat.washington.edu

Department of Statistics  
University of Washington

January 2025

CDF and PDF. Sampling  *Refresher* 

Examples of continuous distributions 

ML estimation for continuous distributions 

ML estimation by gradient ascent

Reading: Ch.5, 6

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

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

Cumulative distribution function (CDF)

← Sampling

$$F(x) = P[X \leq x] \quad (1) \quad \equiv P(-\infty, x]$$

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)

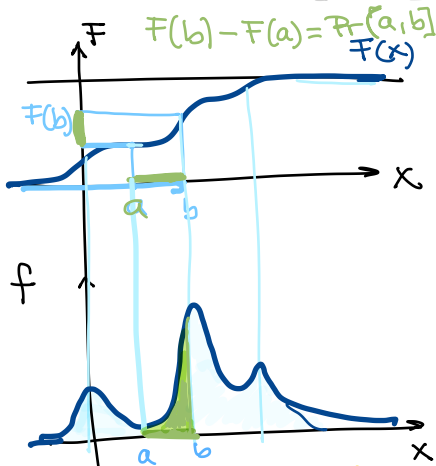
$$f = \frac{dF}{dx} \quad (2)$$

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

Newton

normalization condition

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



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

$\Pr(a, b) = \text{area under } f(x)$

$$\Pr[S] = 1 \Rightarrow \int_{-\infty}^{\infty} f(x) dx = 1$$

$$f(x) \geq 0$$

# 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

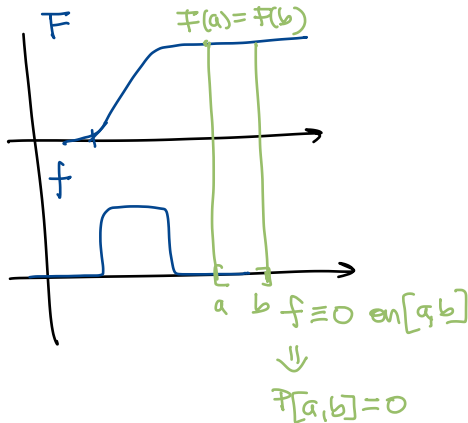
## Probability density [function] (PDF)

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

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

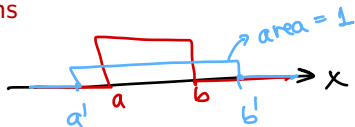
## normalization condition

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



# Examples of continuous distributions

## Model classes



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

uniform

(5)

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

(6)

$\rightarrow a, b = \text{parameters} \in (-\infty, \infty)$   
 $a < b$

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

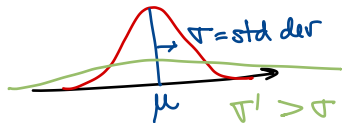
normal

(7)

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

(8)

$\rightarrow \mu, \sigma^2 > 0$



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

(9)

logistic



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

(10)

$\rightarrow a, b$   
 $a > 0$

exponential

$$f(x; \lambda) = \lambda e^{-\lambda x} \quad x \geq 0$$

$\rightarrow \underline{\underline{\lambda > 0}}$

$S = [0, \infty)$



# ML estimation for continuous distributions

## STATISTICS

Given  $x^1, x^2, \dots, x^n \in (-\infty, \infty) = S \sim \text{iid from } P$

what is  ~~$P$~~ ?

$f(x) = ?$

Density estimation

