

STAT 391

4/11/23

Lecture 5

> next Thursday : Q1
HW 1 due tomorrow

Quiz 1:
all material until <
Continuous Sample
Spaces
At beginning of class

Lecture Notes II – Maximum Likelihood Estimation for Discrete Distributions

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

April, 2021

Max Likelihood Principle ✓

ML estimation for arbitrary discrete distributions ✓

Other ML estimation examples ✓

ML estimate as a random variable ←

Reading: Ch. 4.1, 4.2

~~ML estimation for arbitrary discrete distributions~~

A random experiment

Coin toss

$$S = \{0, 1\}$$

$$n=10 \Rightarrow S = \{0, 0.1, 0.2, \dots, 0.5, 1.0\}$$

n_1	θ_1^{ML}	$\theta_1^{\text{true}} = 0.5$	$\theta^{\text{true}} \in \{\theta^{\text{ML}} \text{ observed}\}$
10 trials	4	0.4	
10 trials	7	0.7	
	6	0.6	
	6	0.6	
	6	0.6	

GUESSES

$\theta^{\text{ML}} = 0.55$ possible? **Nb**

- values near θ^{true} are more likely

$$n_1^* = 29 \Rightarrow \theta_1^{*\text{ML}} = 0.58$$

$n^* = 50$ indep trials

take mode $\theta^{\text{mode}} = 0.6$ ↗
 average θ^{ML} / S
 median θ^{ML} 's
 resample ...
 Modern Robust Method

~~ML estimation for arbitrary discrete distributions~~

Candy sampling

θ_{true}
Green

unknown

n	n_{Green}	$\hat{\theta}_{\text{Green}}$
2	0	0
5	1	0.2
2	1	0.5
4	1	0.25
2	1	0.5
6	3	0.5
5	2	0.4

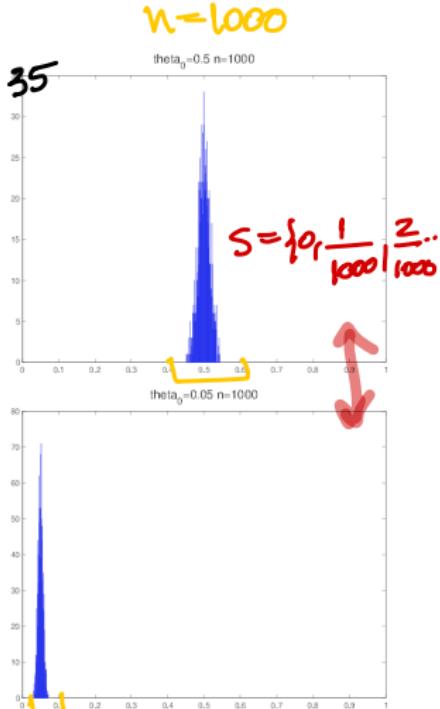
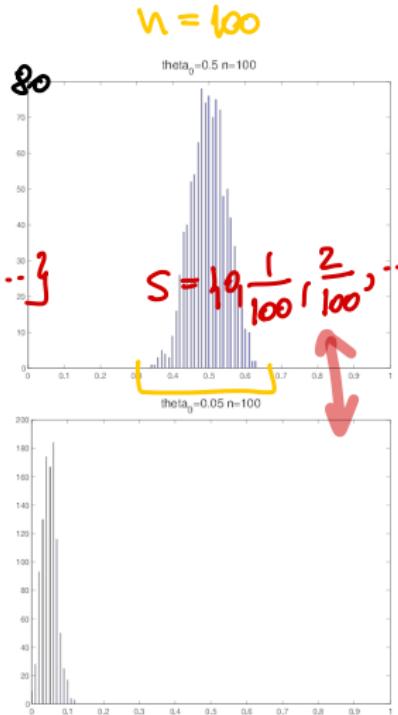
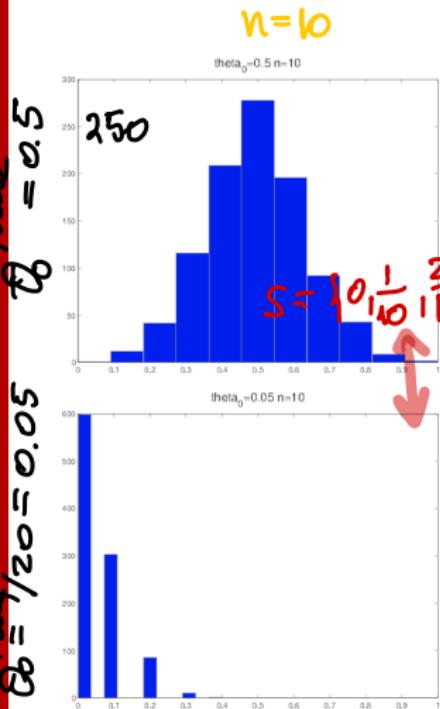
↑
independent

$$n^* = 26$$

$$n_G^* = 9$$

$$\hat{\theta}^{* \text{MC}} = \frac{9}{26} \approx 0.31$$

ML estimate as a random variable $\times 1000$ experiments



distribution of $\hat{\theta}_0^{\text{ML}}$ concentrates \Leftrightarrow

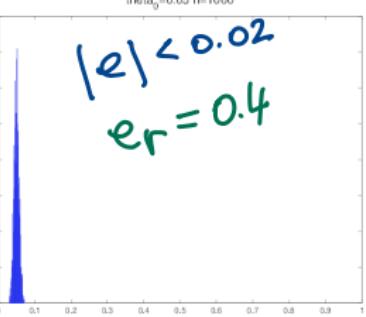
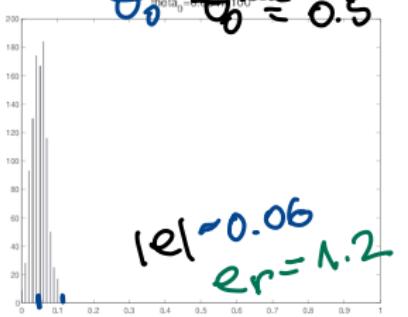
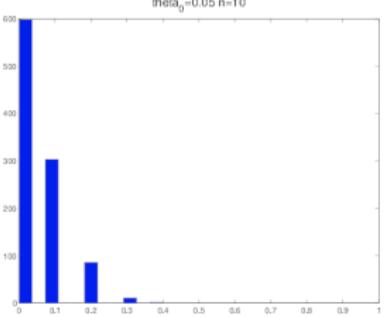
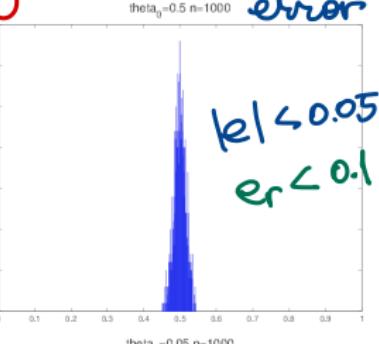
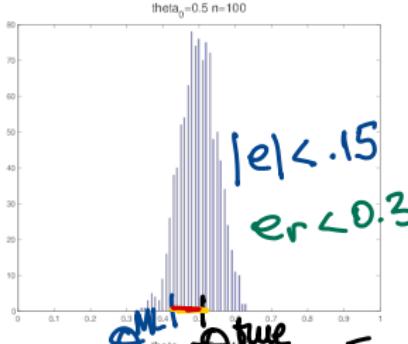
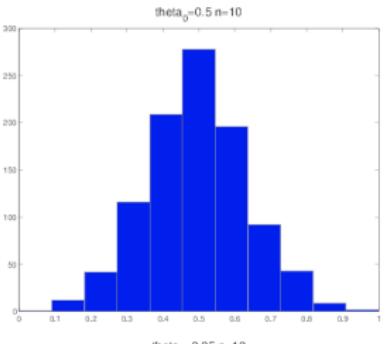
interval that contains observed $\hat{\theta}_0^{\text{ML}}$ shrinks

ML estimate as a random variable

$$\text{Var Bern}(p) = p(1-p)$$

$$e = \theta^{\text{true}} - \theta^{\text{ML}}$$
$$e_r = \frac{|\theta^{\text{true}} - \theta^{\text{ML}}|}{\theta^{\text{true}}}$$

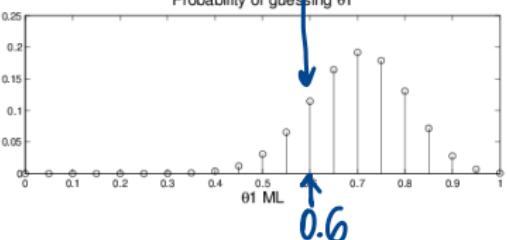
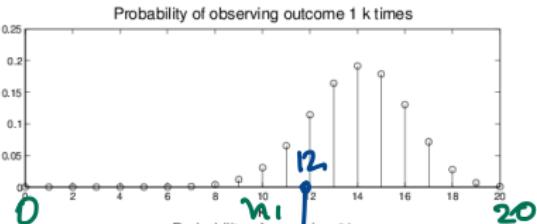
absolute error
relative error



$$\theta^{\text{true}} = 0.05$$

ML estimate as a random variable

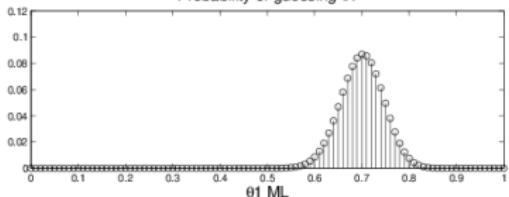
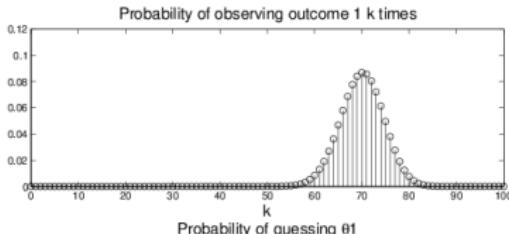
$$n=20 \Rightarrow n_1 \sim \text{Binom}(0.7, 20)$$



$$\hat{\theta}_1^{\text{ML}} = \frac{n_1}{20}$$

$$P\left[\hat{\theta}_1^{\text{ML}} = \frac{n_1}{n} \mid n=20\right] = P[n_1 \mid n=20] = \text{Binom}(\dots)$$

$$\binom{n}{n_0 n_1 \dots n_{m-1}} = \frac{n!}{n_0! n_1! \dots n_{m-1}!}$$



Lecture Notes III: Discrete probability in practice – Small Probabilities

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

April, 2021

The problem with estimating small probabilities

Definitions and setup

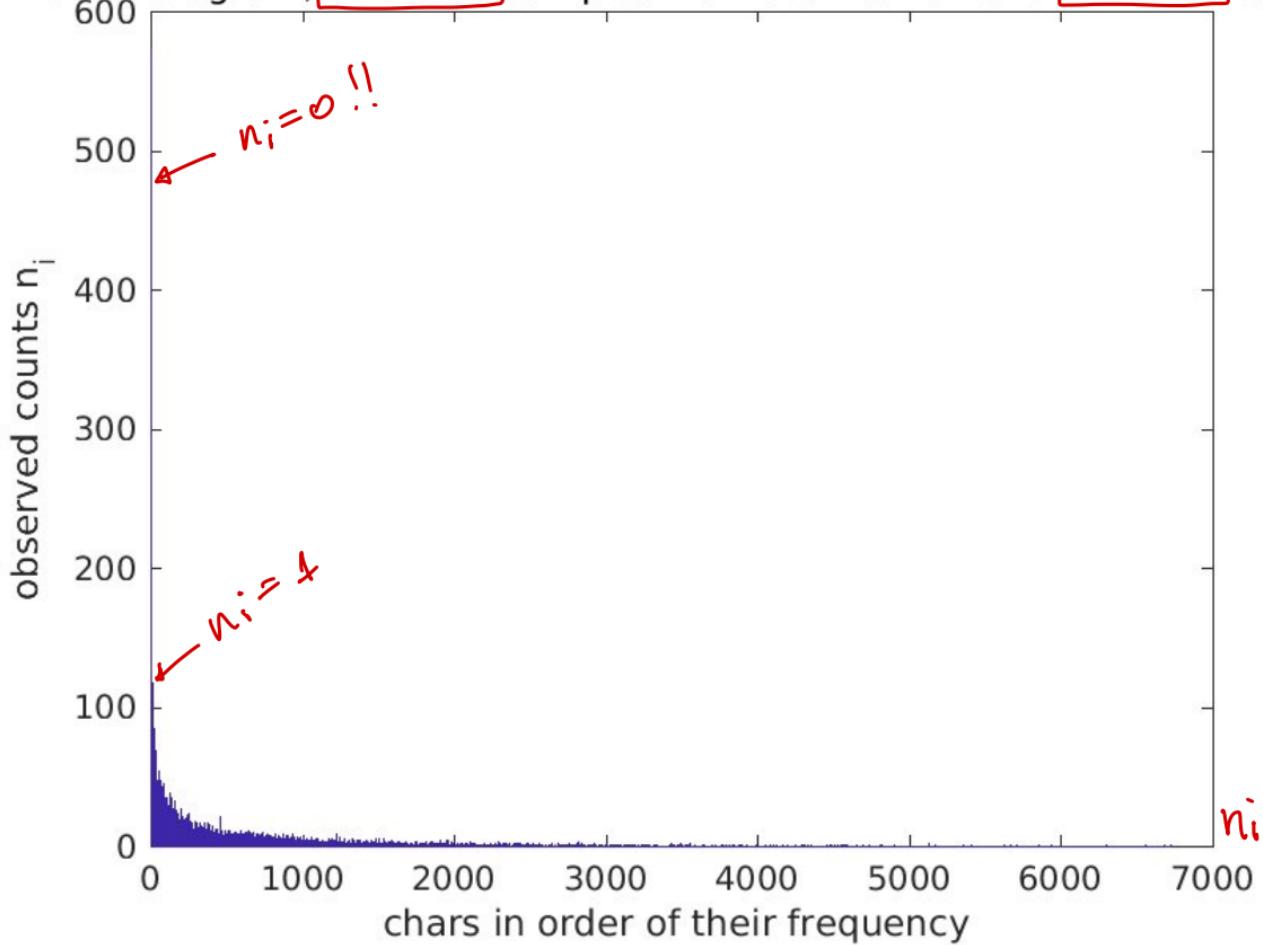
Additive methods (Laplace, Dirichlet, Bayesian, ELE)

Discounting (Ney-Essen)

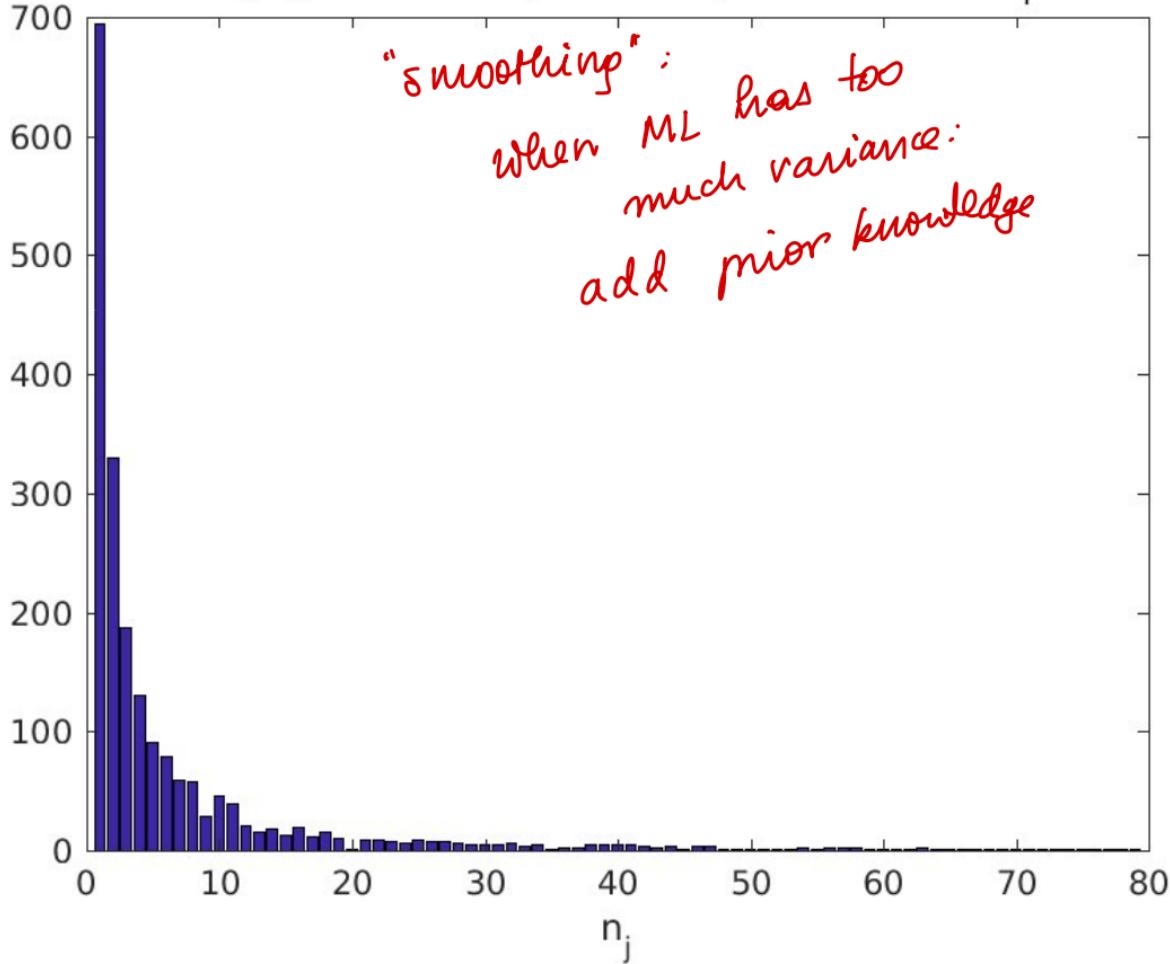
Multiplicative smoothing: Estimating the next outcome (Witten-Bell, Good-Turing)

Back-off or shrinkage – mixing with simpler models

Data histogram, $n=15000$ samples from distribution over $m=9933$ chars



same fingerprint for $k > 0$, $n = 15000$, $m = 9933$ max $n_i = 572$



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

$m = |S|$

i-th letter in sentence
j-th letter in s

- 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 $R_k = \{j, n_j = k\}$

Histogram of hist

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

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

$n_j = 2: e, i$

Histogram of n_j 's

Fingerprint r_k

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

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

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

$r_3 = \dots r_n = 0$

need counts

$$\theta_a^{\text{ML}} = 0 = \theta_b^{\text{ML}}$$

$$\theta_c^{\text{ML}} = \frac{1}{16}$$

BAD!!

NEED SMOOTHING

Principle: if $n_j = n_{j'}$ $\Rightarrow \theta_j = \theta_{j'}$ after smoothing
 $j, j' \in S$