# Analysis of Min Count

## Last time

- $\mathcal{M} = (S, m)$, $S = \{s_1, \ldots s_n\}$, $n = |S| = ?$
- $MinCount : \mathcal{M} \longrightarrow h(s_i) = u_i \sim \mathcal{U}(0,1) \longrightarrow U_1, U_2, \ldots U_n \longrightarrow$

$$\longrightarrow U_{1:n} \leq U_{2:n} \leq \ldots \leq U_{k:n} \longrightarrow \underset{\underset{estimator}{\wedge}}{\hat{n}_k} = \frac{n-1}{U_{1:n}}$$

- $U_{k:n} \sim Beta(k, n+1-k)$, $f_k(x) = \dfrac{x^{k-1}(1-x)^{n-k}}{B(k,n+1-k)}$ for $x \in (0,1)$, where

$$B(k, n+1-k) = \frac{\Gamma(k)\,\Gamma(n+1-k)}{\Gamma(n+1)} \underset{\underset{n,k \in \mathbb{N}}{\uparrow}}{=} \frac{(k-1)!\,(n-k)!}{n!}$$

---

## Estimator Construction

$k=1$ : $U_{1:n} \sim Beta(1,n)$, $f_1(x) = n(1-x)^{n-1}$ $\qquad B(\alpha,\beta) = \int_0^1 t^{\alpha-1}(1-t)^{\beta-1}\,dt$

$$\cdot E[U_{1:n}] = \int_0^1 x \cdot f_1(x)\,dx = n \int_0^1 x(1-x)^{n-1} = n \cdot B(2,n) =$$

$$= n \cdot \frac{\Gamma(2)\,\Gamma(n)}{\Gamma(n+2)} = n \cdot \frac{1!\,(n-1)!}{(n+1)!} = \frac{1}{n+1}$$

then can we use it to estimate $\hat{n} \overset{?}{=} \dfrac{1}{U_{1:n}} \longrightarrow E[\hat{n}] = E\left[\dfrac{1}{U_{1:n}}\right] =$

$$= \int_0^1 \frac{1}{x} f_1(x)\,dx = \int_0^1 \frac{1}{x} n(1-x)^{n-1}\,dx = \infty$$

$\underset{\substack{\uparrow \\ x \text{ is close to } 0 \text{ but the} \\ \text{integral disconverge}}}{}$

generally $E[g(x)] = \int_0^1 g(x) f(x)\,dx =$

$X \sim f(x)$

now let's try for $k \geq 2$

$\bullet \ E[\widehat{u_{K:n}}] = \int_0^1 \frac{1}{x} f_k(x) dx = \int_0^1 \frac{1}{x} \cdot \frac{x^{k-1}(1-x)^{n-k}}{B(k, n+1-k)} dx = \frac{1}{B(k, n+1-k)} \int_0^1 x^{k-2}(1-x)^{n-k} dx$

$\overset{Ex \ 13}{=} \ldots = \frac{n}{k-1}$

$\bullet$ then we can construct estimator $\hat{n}_k := \frac{k-1}{u_{K:n}} \rightarrow E[\hat{n}_k] = n$

---

**Variance**

now we want to calculate variance but

$$Var[X] = \bar{E}[X^2] - E[X]^2$$

$\hookrightarrow$ but then $g(x) = \frac{1}{x^2}$ and integral still doesn't converge. So we do it for $k \geq 3$

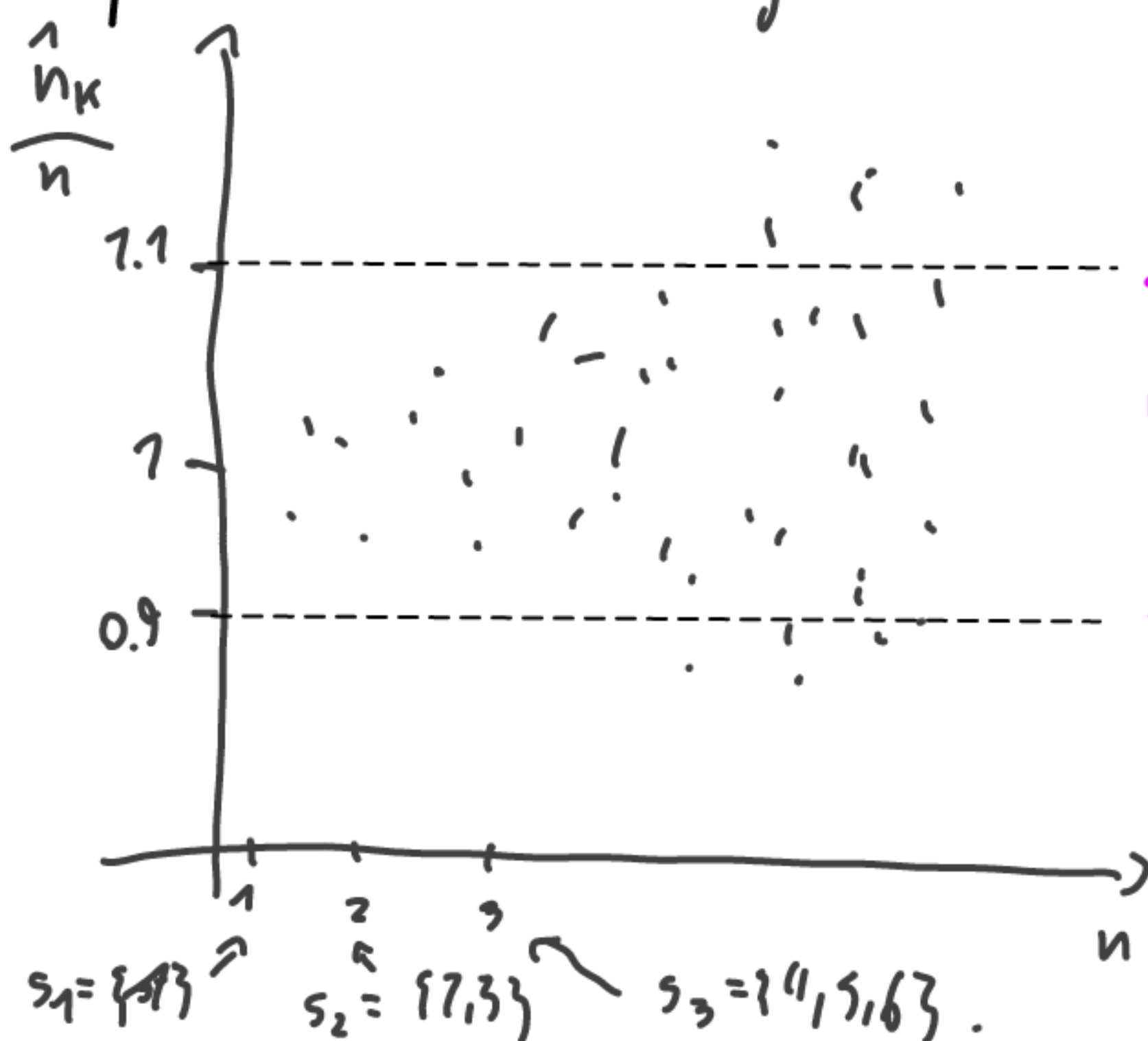$$Var[\hat{n}_k] = \frac{n(n-k+1)}{k-2}$$

---

**Standard Error**

$$SE[\hat{n}_k] = \sqrt{Var\left[\frac{\hat{n}_k}{n}\right]}$$

←  $\color{magenta}\text{this can be interpreted as average error expressed in percente}$

$\approx \frac{1}{\sqrt{k-2}}$ , $k = 100 \rightarrow SE[\hat{n}_k] \approx 10\%$
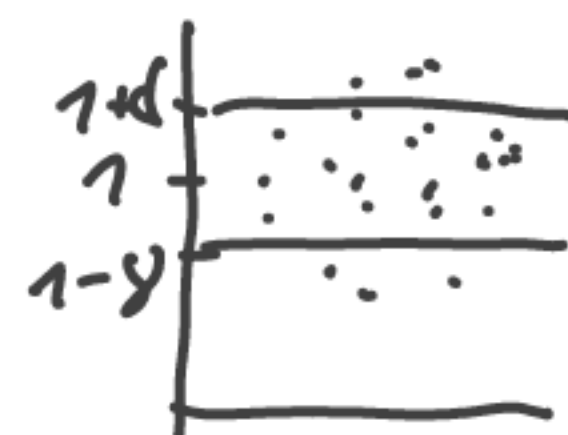
---

**Experiment**    we try to make experiment independent



$\color{magenta}\text{with high probability we would like to say that 99\% of points are in}$

$s_1 = \{9\}$    $s_2 = \{7,3\}$    $s_3 = \{9,5,6\}$ .

SE gives only average error and we would like to know

$$\Pr\left[\;|\frac{\hat{n}_k}{n}-1|\;>\delta\;\right]<\alpha \;\longleftrightarrow\; \Pr\left[\;|\frac{\hat{n}_k}{n}-1|<\delta\;\right]>1-\alpha$$

we will use Chebyshev inequality (Ex. 17)

Let $\delta>0$ and $X$ such that $E[X]$, $Var[X]$ are finite.
Then $\Pr[|X-EX|<\delta]>1-\dfrac{Var[X]}{\delta^2}$

Example $X=\dfrac{\hat{n}_k}{n}$, $E[X]=1$, $\delta=10\%$
$k=100$

$$\Pr\left[\;|\frac{\hat{n}_k}{n}-1|<10\%\right]>1-\dfrac{Var[\frac{\hat{n}_k}{n}]}{(10\%)^2}\approx$$

$$1-\dfrac{\frac{1}{k}}{\frac{1}{100}}=0 \quad\text{, this information is}$$
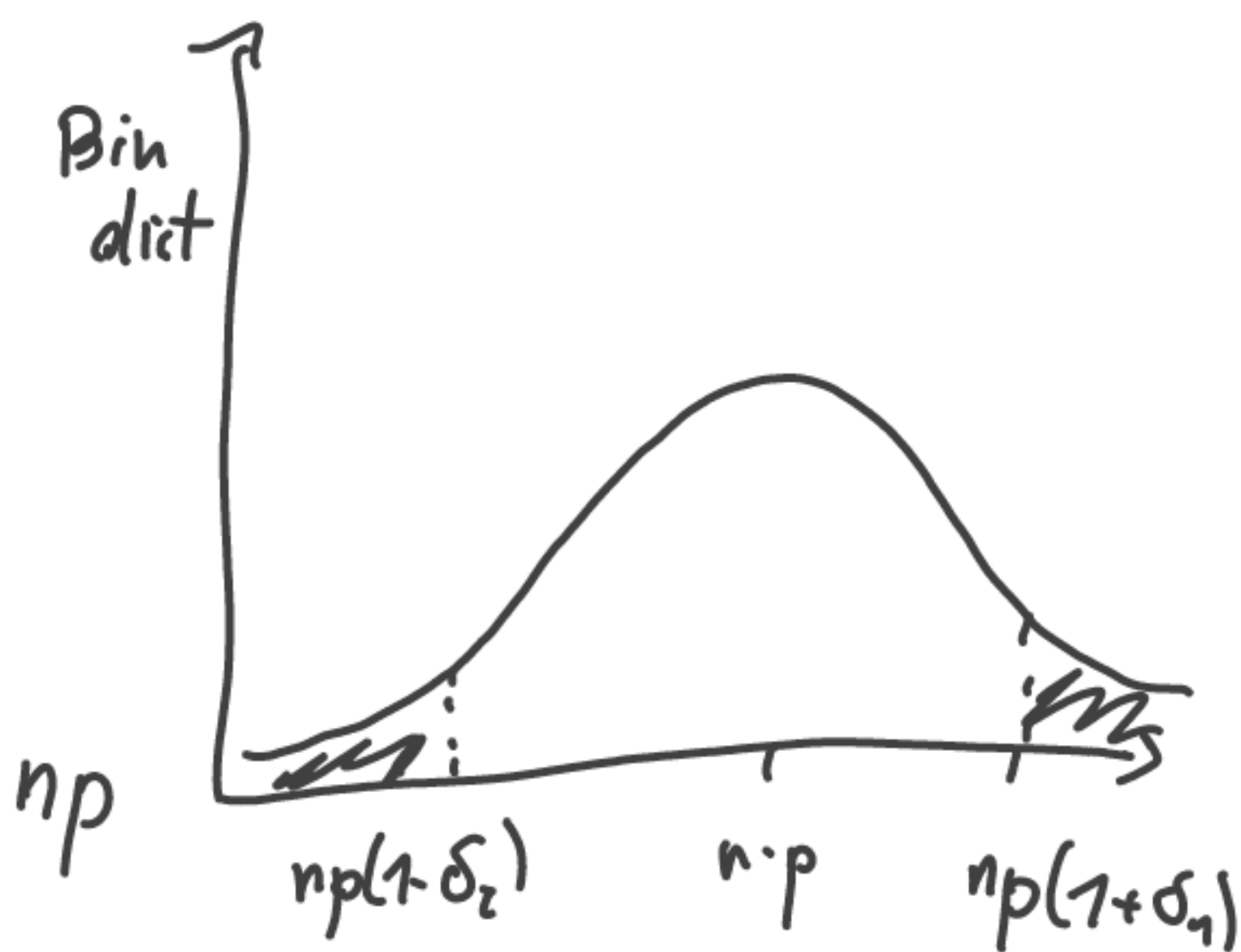
useless.

---

Chernoff inequality (Ex. 18)
Version for binomal distribution $B_{n,p}\sim Bin(n,p)$   $E[B_{n,p}]=n\cdot p$

· $\Pr[B_{n,p}=k]=\binom{n}{k}p^k(1-p)^{n-k}$

· For $\delta_1>0$ and $0<\delta_2<1$ we have

⊛ $\Pr[B_{n,p}\geq n\cdot p(1+\delta_1)]\leq\left(\dfrac{e^{\delta_1}}{(1+\delta_1)^{1+\delta_1}}\right)^{np}$

⊛⊛ $\Pr[B_{n,p}\leq n\cdot p(1-\delta_2)]\leq\left(\dfrac{e^{-\delta_2}}{(1-\delta_2)^{1-\delta_2}}\right)^{np}$

Bin
dist

$np$

$np(1-\delta_2)$    $n\cdot p$    $np(1+\delta_1)$

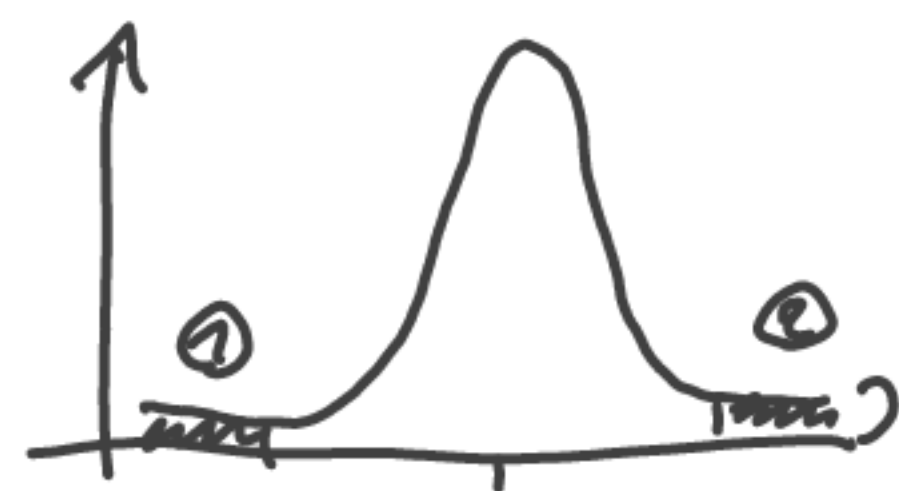# Lemma 1 (Binomial Chernoff → Order Statistic Chernoff)

Let $U_1, U_2, \ldots, U_n$ i.i.d., $U_i \sim U([0,1])$. Let $k \leq n$ and $d \in (0,1)$. Then

$$E[U_{k:n}] = \frac{k}{n}$$

1) if $d < \frac{k}{n}$ then $\Pr[U_{k:n} \leq d] \leq e^{-dn}\left(\frac{dne}{k}\right)^k$

2) if $d > \frac{k}{n}$ then $\Pr[U_{k:n} \geq d] \leq e^{-dn}\left(\frac{dne}{k}\right)^k$

$$E[U_{k:n}] = \frac{k}{n}$$

## Proof:

- $\mathbb{1}_i = \begin{cases} 1 & \text{if } U_i \leq d \\ 0 & \text{if } U_i > d \end{cases}$, $B_{n,d} = \mathbb{1}_1 + \mathbb{1}_2 + \mathbb{1}_3 + \ldots + \mathbb{1}_n \sim \text{Bin}(n,d)$

- note that $U_{k:n} \leq d \iff |\{i : U_i \leq d\}| \geq k \iff B_{n,d} \geq k = dn\left(1 + \frac{k-dn}{dn}\right)$

1) $d < \frac{k}{n}$: $\Pr[U_{k:n} \leq d] = \Pr\left[B_{n,d} \geq \underbrace{dn\left(1 + \frac{k-dn}{dn}\right)}_{k}\right] \overset{\circledast}{\leq}$

$$\leq \left(\frac{e^{\frac{k-dn}{dn}}}{\left(\frac{k}{dn}\right)^{\left(\frac{k}{dn}\right)}}\right)^{nd} = e^{-dn}\left(\frac{dne}{k}\right)^k$$

2) $d > \frac{k}{n}$: $\Pr[U_{k:n} \geq d] \overset{\circledast\circledast}{\leq} e^{-dn}\left(\frac{dne}{k}\right)^k$ $\square$
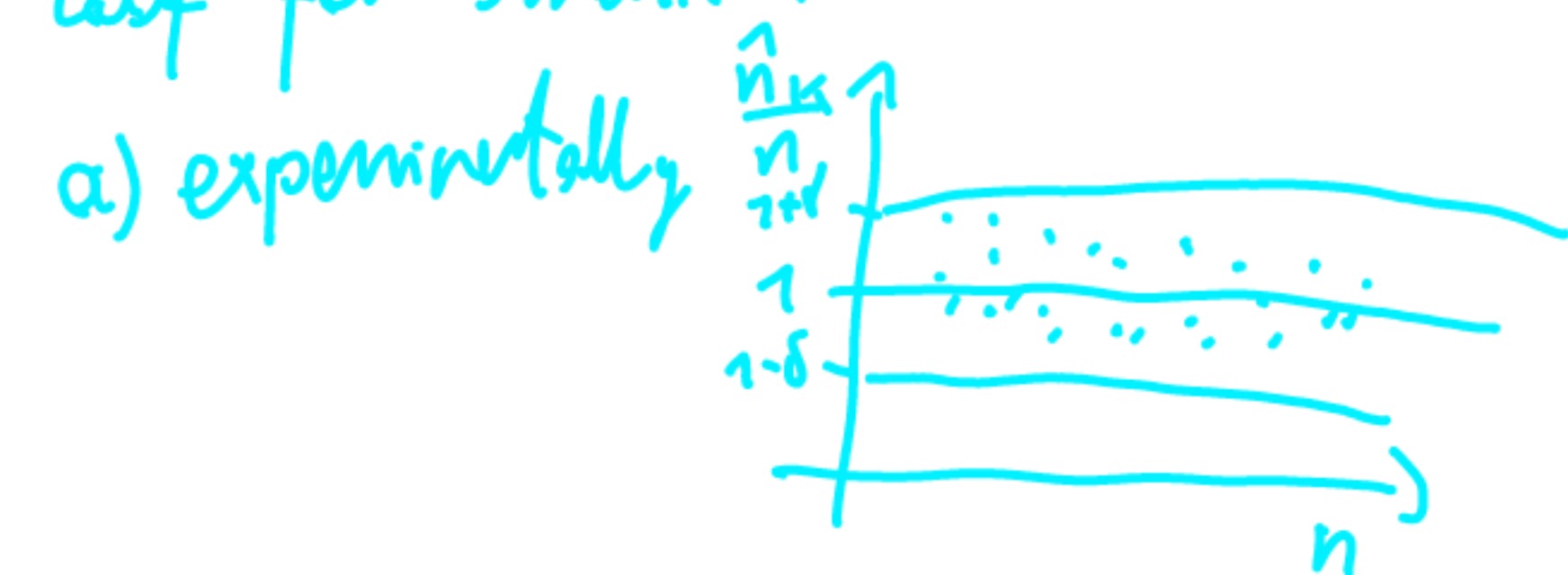
# Theorem 1 (Chernoff bounds for $\hat{n}_k$)

Let $3 \leq k \leq n$, $\varepsilon_1 > 0$, $0 < \varepsilon_2 < 1$

Denote $f_k(x) = e^{xk}(1-x)^k$, then for $\hat{n}_k = \frac{k-1}{u_{k:n}}$ we have

$$Pr\left[\frac{1}{1+\varepsilon_1} \cdot \frac{k-1}{k} < \frac{\hat{n}_k}{n} < \frac{1}{1-\varepsilon_2} \cdot \frac{k-1}{k}\right] > 1 - f_k(\varepsilon_2)^{\frac{\delta}{1+\delta}} - f_k(-\varepsilon_1)^{\frac{\delta}{1-\delta}}$$

$\approx 1$

$(1-\delta)$     $\uparrow$     $(1+\delta)$    $\approx 1$

$\varepsilon_1 = \frac{\delta}{1-\delta_1}$      $\varepsilon_2 = \frac{\delta}{1+\delta}$

---

notes for Task 7

Look for smallest $\delta$ such that $Pr\left[1-\delta < \frac{\hat{n}_k}{n} < 1+\delta\right] > 1-d$, $k=400$, $d=1\%$

a) experimentally



b) use Chebyshev inequality

c) use Chernoff (Theorem 1)

---

Proof: it's obvious

① $Pr[a < \mathcal{X} < b] \stackrel{A}{=} Pr[\mathcal{X} < b] - Pr[\mathcal{X} < a]$



then $Pr\left[\frac{1}{1+\varepsilon_1} \frac{k-1}{k} < \frac{\hat{n}_k}{n} < \frac{1}{1-\varepsilon_2} \frac{k-1}{n}\right] \stackrel{A}{=} Pr\left[\frac{\hat{n}_k}{n} < \frac{1}{1-\varepsilon_2} \frac{k-1}{n}\right] - Pr\left[\frac{\hat{n}_k}{n} < \frac{1}{1+\varepsilon_1} \frac{k-1}{k}\right]$

$\underbrace{}_{\stackrel{\geq}{③} 1 - f_k(\varepsilon_2)}$      $\underbrace{}_{\stackrel{\leq}{②} f_k(-\varepsilon_1)}$

② $Pr\left[\frac{\hat{n}_k}{n} < \frac{1}{1+\varepsilon_1} \frac{k-1}{k}\right] = Pr\left[\frac{k-1}{u_{k:n}} < \frac{n}{1+\varepsilon_1} \frac{k-1}{k}\right] = Pr\left[u_{k:n} > \underbrace{\frac{k}{n}(1+\varepsilon_1)}_{a}\right] \stackrel{\text{Lemma 1}}{\leq} e^{-dn\left(\frac{dn+e}{k}\right)^k}$

$= \ldots =$

$= f_k(-\varepsilon_1)$

③ similar reasoning as ② $P_n \left[ \frac{\hat{n}_n}{n} < \frac{k-1}{k} \frac{1}{1+\xi_2} \right] \geq 1 - f_k(\xi_2) \Longleftrightarrow$

$$- \quad | \quad | \quad - \quad \leq f_k(\xi_2) \quad \square$$