

Convergence

钟宛君 (18110955) , 黄羽盼 (18110943)

Nov 1

Content

Part A: Sampling and Estimation

Part B: Probably Approximately Correct (PAC)

Part C: Concentration of Measure: Markov, Chebyshev Inequality and Chernoff-Hoeffding Inequality

Content

Part A: Sampling and Estimation

Part B: Probably Approximately Correct (PAC)

Part C: Concentration of Measure: Markov, Chebyshev Inequality and Chernoff-Hoeffding Inequality

Sampling and Estimation

- Data: Set of n data points $\mathbf{P} = \{p_1, p_2, \dots, p_n\}$
- Powerful Assumption: Data comes iid (**Identically** and **Independent Distributed**) from **fixed, unknown** PDF (Probability density function).
- Our goal is to **estimate the underlying data distribution**.

Sampling and Estimation

➤ What's the mean of PDF?

Considering random variable $X : X \sim f$:

mean of f is the expected value of X : $\mathbf{E}[X]$

➤ How we estimate the mean of f ? - By sample mean

$$\bar{P} = \frac{1}{n} \sum_{i=1}^n p_i$$

We can randomly sample n data points to estimate the mean of f by sample mean.

$$\bar{P} = \frac{1}{n} \sum \{p_i\} \leftarrow \text{realize } \{X_i\} \underset{\text{iid}}{\sim} f$$

Step 1: Select n iid variables $\{X_i\}$ corresponding to set of n independent observation $\{p_i\}$

Step 2: Take their average to estimate the mean of f .

Central Limit Theorem

- Goal: Estimate how well the sample mean approximates the true mean
- The sample mean is dependent to the data we select ($\{X_i\}$ is randomly selected). Therefore it's not precisely equal to the mean of f .
- Central Limit Theorem

Central Limit Theorem: Consider n iid random variables X_1, X_2, \dots, X_n , where each $X_i \sim f$ for any fixed distribution f with mean μ and bounded variance σ^2 . Then $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ converges to the normal distribution with mean $\mu = \mathbf{E}[X_i]$ and variance σ^2/n .

Which means: The mean of random sample mean will converge to the normal distribution with mean $\mathbf{E}[\mathbf{X}]$ as the observation increases.

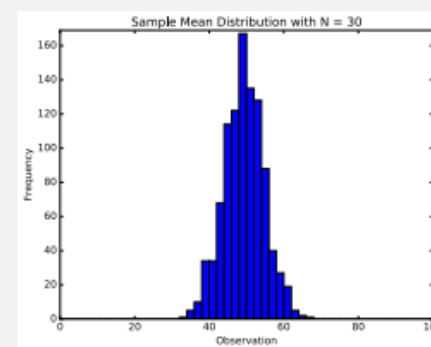
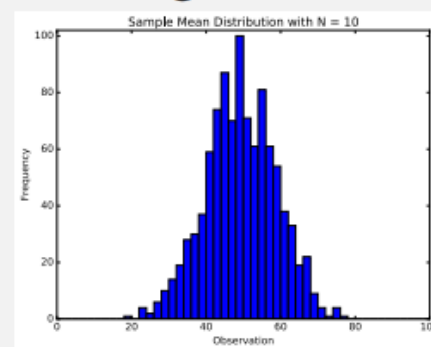
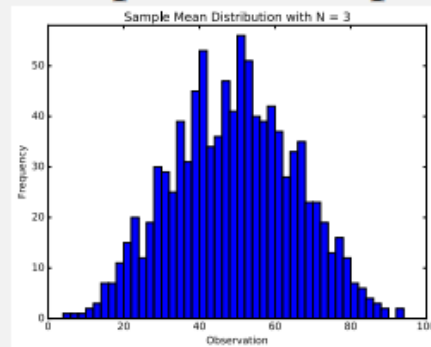
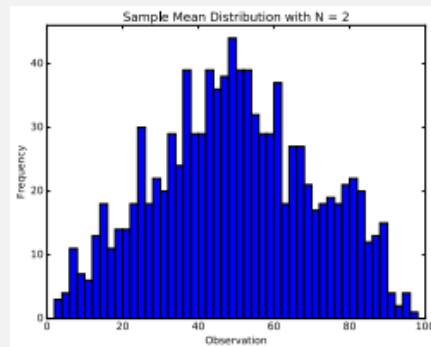
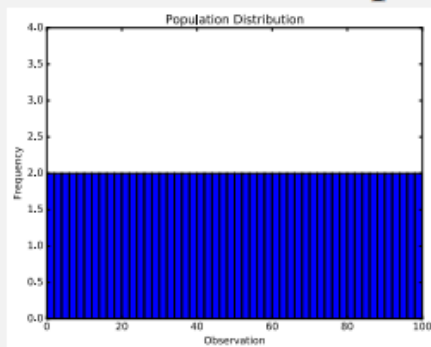
$\{p_i\}$

Central Limit Theorem

➤ Example

Example: Central Limit Theorem

Consider f as a uniform distribution over $[0, 100]$. If we create n samples $\{p_1, \dots, p_n\}$ and their mean \bar{P} , then repeat this 1000 times, we can plot the output in histograms:



We see that starting at $n = 2$, the distributions look vaguely normal (in the technical sense of a normal distribution), and that their standard deviations narrow as n increases.

Central Limit Theorem Example

Rolling n dice, each governed by the same probability distribution. Each face of a dice has probability $1/6$.

mean

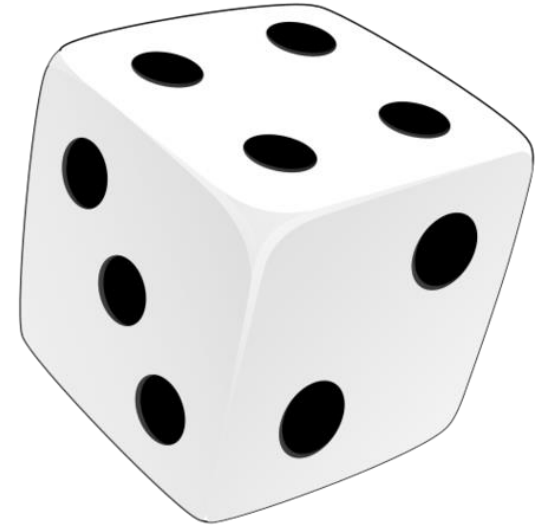
$$\mu = E(x_i) = \frac{1 + 2 + 3 + 4 + 5 + 6}{6} = 3.5$$

variance

$$\frac{\sigma^2}{n} = \frac{(2.5^2 + 1.5^2 + 0.5^2) * 2}{6} \approx 2.92$$

standard deviation

$$\sqrt{\frac{\sigma^2}{n}} \approx 1.71$$



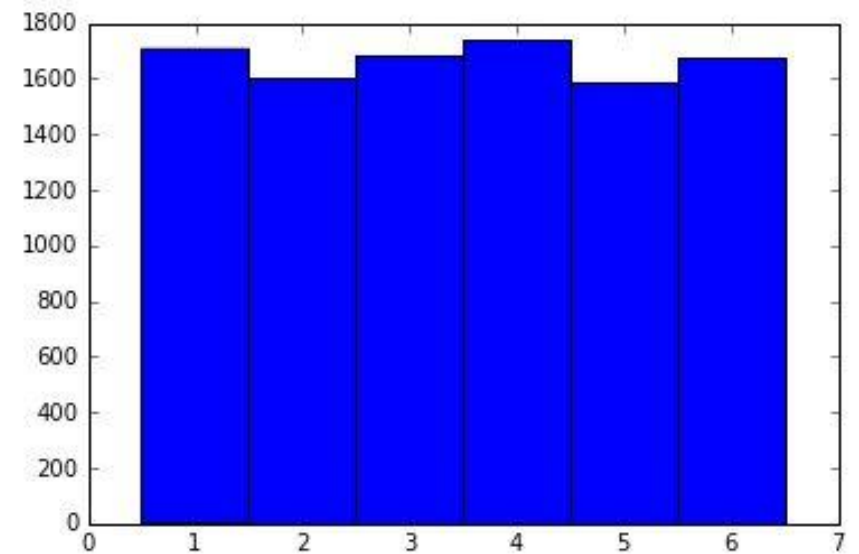
Reference: <https://zhuanlan.zhihu.com/p/25241653> and <https://www.albany.edu/~jr853689/CentralLimitTheoremForDice.htm>

You can do Dice Experiment on <http://www.randomservices.org/random/apps/DiceExperiment.html>

Central Limit Theorem Example

1. Rolling a dice for 10000 times.
Draw the histogram of random data.

```
>>> import numpy as np
>>> random_data = np.random.randint(1, 7, 10000)
>>> random_data.mean()
'3.5199'
>>> random_data.std()
'1.7116085971973851'
```



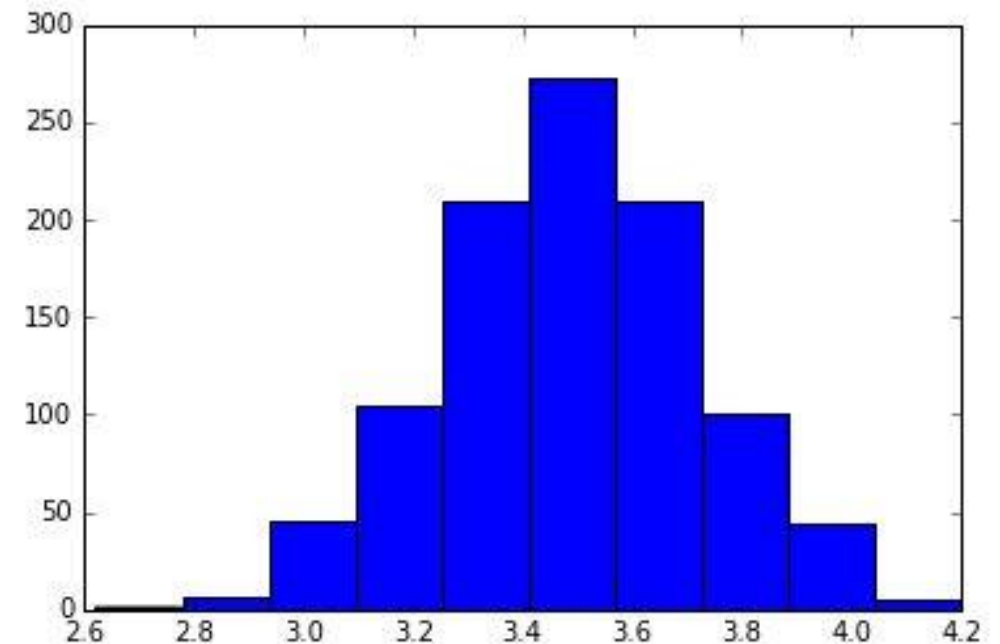
2. Choose ten scores from the generated data randomly and calculate their mean and variance.

```
[>>> sample1 = []
>>> for i in range(0, 10):
[...     sample1.append(random_data[int(np.random.random() * len(random_data))])
[...
>>> np.mean(sample1)
'2.5'
>>> np.std(sample1)
'1.8027756377319946'
```

3. Now choose 1000 groups, each group has 50 samples. Draw the histogram of the 1000 mean scores.

```
>>> samples_mean = []
>>> samples_std = []
>>>
>>> for i in range(0, 1000):
...     sample = []
...     for j in range(0, 50):
[...         sample.append(random_data[int(np.random.random() * len(random_data))])
[...         samples_mean.append(np.mean(sample))
[...         samples_std.append(np.std(sample))
[...
[>>> np.mean(samples_mean)
'3.51656'
[>>> np.mean(samples_std)
'1.6926817755030283'
```

- population average: 3.5 is the "average of all possible rolls of a fair die."
- The output and data distribution illustrate the Central Limit Theorem



Remaining Mysteries

➤ What does convergence mean?

Convergence refers to **what happens as some parameter increases**. (eg. n goes to infinity then the distribution will be more precise.)

➤ How we formalize the error when estimating?

The distance between $\bar{\mathbf{P}}$ and μ is more than ϵ , with probability at most δ .

We call that **"probably approximately correct"** (PAC).

➤ How we describe the distribution of f more detailedly?

We discuss some very common concentration of measure tools - To **provide the upper bounds** to state the PAC bounds.

Content

Part A: Sampling and Estimation

Part B: Probably Approximately Correct (PAC)

Part C: Concentration of Measure : Markov, Chebyshev Inequality and Chernoff-Hoeffding Inequality

Probably Approximately Correct (PAC)

➤ Concentration of measure bounds

Introduce three most common concentration of measure bounds, provide increasingly strong bounds, but requires increasing information about underlying f

- a. Markov Inequality
- b. Chebyshev Inequality
- c. Chernoff-Hoeffding Inequality

➤ Basic form of PAC bound

$$\mathbf{Pr}[|\bar{X} - \mathbf{E}[\bar{X}]| \geq \varepsilon] \leq \delta.$$

Content

Part A: Sampling and Estimation

Part B: Probably Approximately Correct (PAC)

Part C: Concentration of Measure : Markov, Chebyshev Inequality and Chernoff-Hoeffding Inequality

Markov Inequality

➤ Theorem

Let X be a non-zero, random variable, and $a > 0$ then

$$\Pr[|X| \geq a] \leq \frac{E[|X|]}{a}$$

Or equivalently $\Pr[|X| \geq aE[|X|]] \leq \frac{1}{a}$.

PAC bound with $\epsilon = a - E[|X|]$ and $\delta = E[|X|]/a$. Then the bound can be rephrased as:

$$\Pr[|X - E[|X|]| \geq \epsilon] \leq \frac{E[|X|]}{\epsilon + E[|X|]}$$

Markov Inequality

➤ **Conclusion:** It provide **weak bounds** but it **only requires only expected value $E[|X|]$ of f**

➤ **Example**

For the toss of n fair coins let X_i denote the event that the i^{th} coin lands heads. Then $E[\sum X_i] = \frac{n}{2}$ so the probability that more than $2/3$'s of the coins come up heads is

$$\Pr \left[\frac{2n}{3} \text{ come up heads} \right] \leq \frac{n/2}{2n/3} = \frac{3}{4}$$

Chebyshev's Inequality

➤ Theorem

$$\Pr[|X - E[X]| \geq \epsilon] \leq \frac{\text{Var}[X]}{\epsilon^2}$$

where $\text{Var}[X] = E[(X - E[X])^2]$. The bound is $\delta = \frac{\text{Var}[X]}{\epsilon^2}$.

- This bound is **typically stronger than the Markov** because δ decreases quadratically in ϵ instead of linearly.
- **Conclusion**: Compared with Markov inequality, it has two property:
 1. Have **stronger PAC bounds**, because δ decreases quadratically in ϵ .
 2. **Support negative value of X** by square function.

Chebyshev's Inequality

➤ **Proof** Let random variable $Y = |X - E[X]|$. Using Markov's inequality we have that

$$\Pr[Y \geq a] = \Pr[Y^2 \geq a^2] \leq \frac{E[Y^2]}{a^2} = \frac{\text{Var}[X]}{a^2}$$

➤ Example

Example: Chebyshev for IID Samples

Recall that for an average of random variables $\bar{X} = (X_1 + X_2 + \dots + X_n)/n$, where the X_i s are iid, and have variance σ^2 , then $\text{Var}[\bar{X}] = \sigma^2/n$. Hence

$$\Pr[|\bar{X} - E[X_i]| \geq \varepsilon] \leq \frac{\sigma^2}{n\varepsilon^2}.$$

Consider now that we have input parameters ε and δ , our desired error tolerance and probability of failure. If can draw $X_i \sim f$ (iid) for an unknown f (with known expected value and variance σ), then we can solve for how large n needs to be: $n = \sigma^2/(\varepsilon^2\delta)$.

Chernoff/Hoeffding Inequality

➤ Theorem

Let X_1, \dots, X_n be independent random variables in the interval $[0, 1]$ and let $X = \sum X_i$.
Then

$$\Pr[|X - E[X]| \geq t] \leq 2\exp(-2t^2/n)$$

Equivalently

$$\Pr[|X - E[X]| \geq \epsilon E[X]] \leq 2\exp(-2\epsilon^2 E[X]^2/n)$$

and

$$\Pr[|X - E[X]| \geq \epsilon n] \leq 2\exp(-2\epsilon^2 n)$$

➤ The bound is

$$\delta = 2\exp(-2\epsilon^2 n)$$

For desired error tolerance ϵ and failure probability δ , we can set

$$n = \left(\frac{1}{2\epsilon^2}\right) \ln\left(\frac{2}{\delta}\right)$$

Although this has similar relationship with ϵ and δ , but the dependence of n on δ is exponentially less for this bound.

Chernoff-Hoeffding Inequality

- Consider set of iid random variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ where $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Now we assume each \mathbf{X}_i lies in a bounded domain $[b, t]$.

$$\Pr[|\bar{X} - \mathbf{E}[\bar{X}]| > \varepsilon] \leq 2 \exp\left(\frac{-2\varepsilon^2 n}{\Delta^2}\right)$$

The bound is

$$\delta = 2 \exp\left(\frac{-2\varepsilon^2 n}{\Delta^2}\right)$$

For desired error tolerance ε and failure probability δ , we can set

$$n = (\Delta^2 / (2\varepsilon^2)) \ln(2 / \delta)$$

Although this has similar relationship with ε and δ , but the dependence of n on δ is exponentially less for this bound.

The Derivation of Three Basic Inequality

Suppose that \mathbf{Z} has a finite mean and that $\mathbb{P}(\mathbf{Z} \geq 0) = 1$. Then, for any $\epsilon > 0$,

$$\mathbb{E}(\mathbf{Z}) = \int_0^{\infty} z dP(z) \geq \int_{\epsilon}^{\infty} z dP(z) \geq \epsilon \int_{\epsilon}^{\infty} dP(z) = \epsilon \mathbb{P}(\mathbf{Z} > \epsilon)$$

which yields *Markov's inequality*:

$$\mathbb{P}(\mathbf{Z} > \epsilon) \leq \frac{\mathbb{E}(\mathbf{Z})}{\epsilon}$$

An immediate consequence of Markov's inequality is *Chebyshev's inequality*

$$\mathbb{P}(|\mathbf{Z} - \mu| > \epsilon) = \mathbb{P}(|\mathbf{Z} - \mu|^2 > \epsilon^2) \leq \frac{\mathbb{E}(\mathbf{Z} - \mu)^2}{\epsilon^2} = \frac{\sigma^2}{\epsilon^2}$$

Where $\mu = \mathbb{E}(\mathbf{Z})$ and $\sigma^2 = \text{Var}(\mathbf{Z})$.

The Derivation of Three Basic Inequality

If Z_1, \dots, Z_n are iid with mean μ and variance σ^2 then, since $\text{Var}(\overline{Z}_n) = \frac{\sigma^2}{n}$, Chebyshev's inequality yields

$$\mathbb{P}(|\overline{Z}_n - \mu| > \epsilon) \leq \frac{\sigma^2}{n\epsilon^2}$$

While this inequality is useful, it does not decay exponentially fast as n increases.

To improve the inequality, we use *Chernoff's method*: for any $t > 0$,

$$\mathbb{P}(Z > \epsilon) = \mathbb{P}(e^Z > e^\epsilon) = \mathbb{P}(e^{tZ} > e^{t\epsilon}) \leq e^{-t\epsilon} \mathbb{E}(e^{tZ})$$

Chernoff-Hoeffding Inequality If Z_1, \dots, Z_n are independent with $\mathbb{P}(a \leq Z_i \leq$

Conclusion

| | Markov | ChebyShev | Chernoff-Hoeffding |
|-----------------------|------------------|----------------------------|------------------------|
| Strength of PAC bound | Weak (Linealy) | Meduim (Quadratically) | Strong (Exponentially) |
| Information requires | $E[x]$ | $E[x]$ and $\text{Var}[x]$ | Bound of X |
| Value of X | Larger than zero | Arbitrary | Arbitrary |

Application 1: Approximating the fraction of 1's in a binary string

Suppose we want to estimate the fraction of 1's in a given string $S \in \{0, 1\}^n$.

That is, we wish to find a fast randomized algorithm that, given ϵ and string S outputs a value V such that $|V - \text{fraction of 1's}| < \epsilon$ with probability $2/3$.

Algorithm. Pick $k = 1/\epsilon^2$ uniformly random indices in the string S and output the fraction of 1's in the sample.

Analysis. Let X_1, \dots, X_k be random variables indicating if a 1 was found in the string position for the i^{th} index selected ($1 \leq i \leq k$). Then by Chernoff's bound

$$\Pr \left[\left| \frac{X}{k} - \frac{E[X]}{k} \right| > \epsilon \right] \leq 2e^{-2\epsilon^2 k} = 2e^{-2} < \frac{1}{3}$$

as $\epsilon^2 k = 1$.

meaning that we output a good estimate (i.e. within ϵ from the true fraction of 1's in the string) with probability $> 2/3$.

Application 2: Improving a random algorithm's correctness

Suppose we are given a randomized algorithm A which on each input x from some domain D outputs a 0 or 1 answer and it is correct with probability $p = 2/3$. Let algorithm B run A for t times and output the majority answer.

Show that algorithm B is correct (on each input) with probability greater than $1 - 2^{-ct}$ for some constant c (that is, $\forall x \in D, \Pr[B(x) = f(x)] \geq 1 - 2^{-ct}$.)

Analysis. Let X_1, \dots, X_t be indicator variables such that $X_i = 1$ if A outputs the correct answer in the i^{th} step. Therefore, $E[X_i] = p = 2/3$. Set $X = \sum X_i$, that is X is the random variable counting the number of correct answers, and notice that $E[X] = 2t/3$.

$$\begin{aligned} & \Pr[B \text{ outputs incorrect answer}] \\ &= \Pr \left[A \text{ outputs incorrect answer more than } \frac{t}{2} \text{ times} \right] = \Pr \left[X < \frac{t}{2} \right] \\ &\leq \Pr \left[X - \frac{2t}{3} < \frac{t}{2} - \frac{2t}{3} \right] \\ &= \Pr \left[X - \frac{2t}{3} < -\frac{t}{6} \right] \leq \Pr \left[\left| X - \frac{2t}{3} \right| > \frac{t}{6} \right] \\ &\leq 2e^{-\frac{2t^2 \left(\frac{1}{6}\right)^2}{t}} = 2^{-ct} \end{aligned}$$

Example

Example: Uniform Distribution

Consider a random variable $X \sim f$ where $f(x) = \{\frac{1}{2} \text{ if } x \in [0, 2] \text{ and } 0 \text{ otherwise.}\}$, i.e, the Uniform distribution on $[0, 2]$. We know $\mathbf{E}[X] = 1$ and $\mathbf{Var}[X] = \frac{1}{3}$.

- Using the Markov Inequality, we can say $\mathbf{Pr}[X > 1.5] \leq 1/(1.5) \approx 0.6666$ and $\mathbf{Pr}[X > 3] \leq 1/3 \approx 0.33333$.
or $\mathbf{Pr}[X - \mu > 0.5] \leq \frac{2}{3}$ and $\mathbf{Pr}[X - \mu > 2] \leq \frac{1}{3}$.
- Using the Chebyshev Inequality, we can say that $\mathbf{Pr}[|X - \mu| > 0.5] \leq (1/3)/0.5^2 = \frac{4}{3}$ (which is meaningless). But $\mathbf{Pr}[|X - \mu| > 2] \leq (1/3)/(2^2) = \frac{1}{12} \approx 0.08333$.

Now consider a set of $n = 100$ random variables X_1, X_2, \dots, X_n all drawn iid from the same pdf f as above. Now we can examine the random variable $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. We know that $\mu_n = \mathbf{E}[\bar{X}] = \mu$ and that $\sigma_n^2 = \mathbf{Var}[\bar{X}] = \sigma^2/n = 1/(3n) = 1/300$.

- Using the Chebyshev Inequality, we can say that $\mathbf{Pr}[|\bar{X} - \mu| > 0.5] \leq \sigma_n^2/(0.5)^2 = \frac{1}{75} \approx 0.01333$, and $\mathbf{Pr}[|\bar{X} - \mu| > 2] \leq \sigma_n^2/2^2 = \frac{1}{1200} \approx 0.0008333$.
- Using the Chernoff-Hoeffding bound, we can say that $\mathbf{Pr}[|\bar{X} - \mu| > 0.5] \leq 2 \exp(-2(0.5)^2 n / \Delta^2) = 2 \exp(-100/8) \approx 0.0000074533$, and $\mathbf{Pr}[|\bar{X} - \mu| > 2] \leq 2 \exp(-2(2)^2 n / \Delta^2) = 2 \exp(-200) \approx 2.76 \cdot 10^{-87}$.

Thanks!

Wanjun Zhong, Yupan Huang