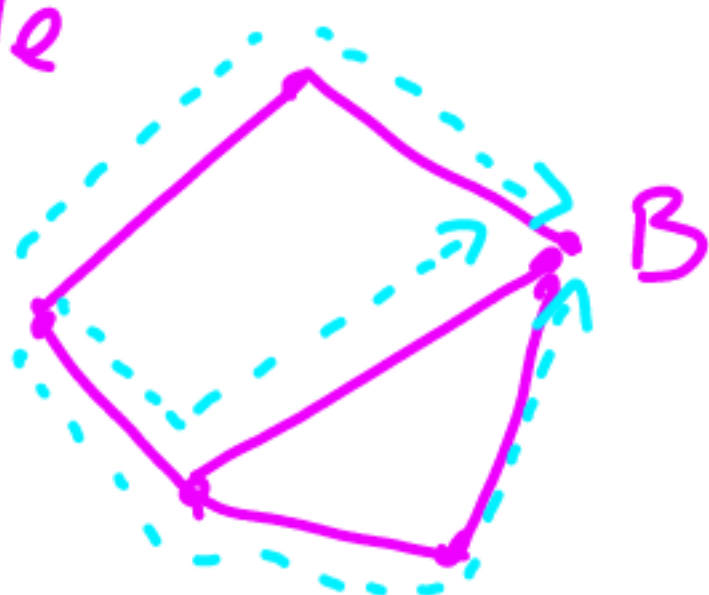


Approximate summation, operations on data sketches

Example

A
(i, λ_i)

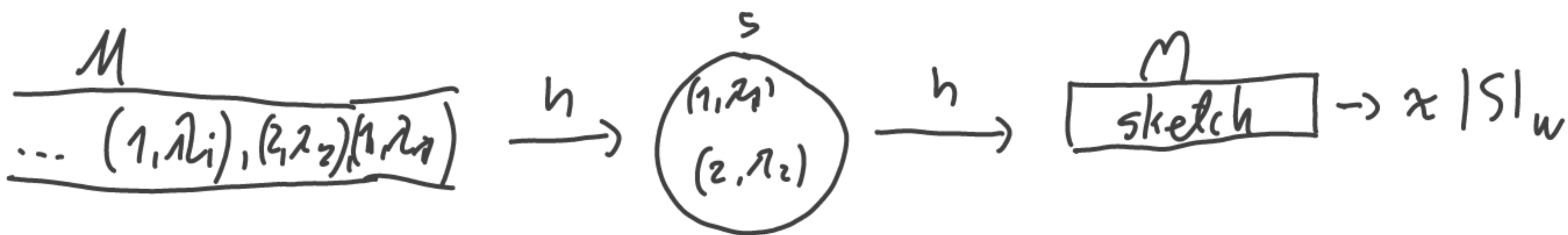


one node (B) can observe one packet many times.

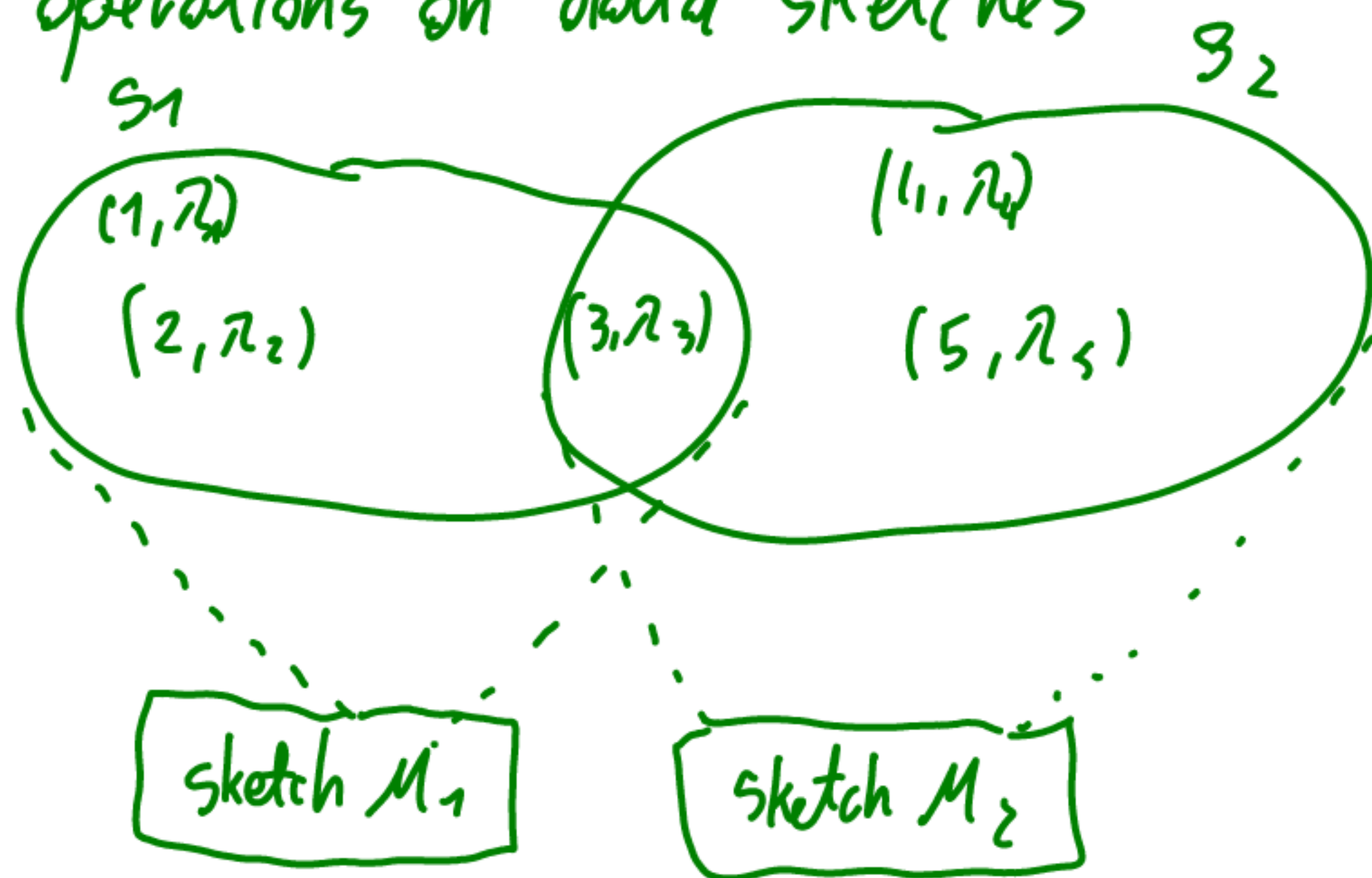
- each node observes a stream $M = (S, m)$, $m: S \rightarrow \mathbb{N}_+$
- $S = \{(1, \lambda_1), (2, \lambda_2), \dots, (n, \lambda_n)\}$, where
 - for (i, λ_i) , i - denotes a unique identifier of a packet
 - $\lambda_i \in \mathbb{R}_+$ denotes size of a packet.

Aim 1: for each node we try to estimate

$$|S|_w := \sum_{(i, \lambda_i) \in S} \lambda_i$$



Aim 2: operations on data sketches



- Use sketches M_1 and M_2 to estimate the value of $|S_1 \cap S_2|_w$, $|S_1 \cup S_2|_w$, $|S_1 \setminus S_2|_w$ very nice

Exp Sketch (M, m, h) :

$$M = (\infty, \infty, \dots, \infty)$$

for each $(i, \lambda_i) \in M$

for each $k \in \{1, 2, \dots, m\}$

$$u \leftarrow h(i \| k)$$

$$M_k \leftarrow \min \left\{ M_k, \frac{\ln(u)}{-\lambda_i} \right\}$$

// concat of binary repr.
we treat $h(i \| k) \sim U(0, 1)$

return M

$$S_i^{(k)} \sim \text{Exp}(\lambda_i)$$

Lemma 1 Let $F(x)$ be CDF of some distribution and let

$$F^{-1}(u) = \inf \{x : F(x) \geq u\}, \quad 0 \leq u < 1.$$

smallest x so that $F(x) \geq u$

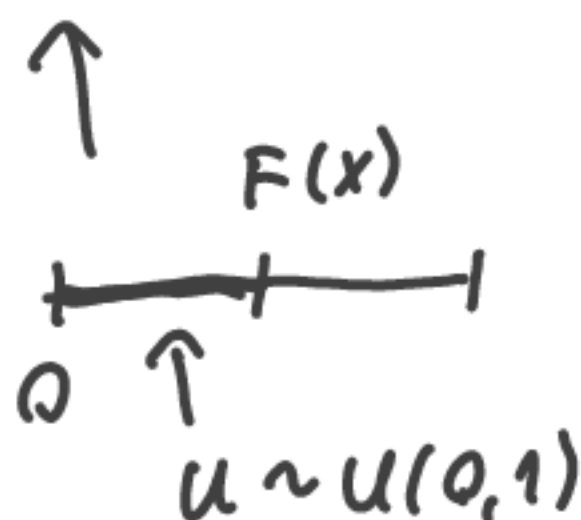
$$\text{CDF} \\ F_X(x) = \Pr[X \leq x]$$

then for $U \sim U(0, 1)$ we have $F^{-1}(U) \sim F(x)$

Proof:

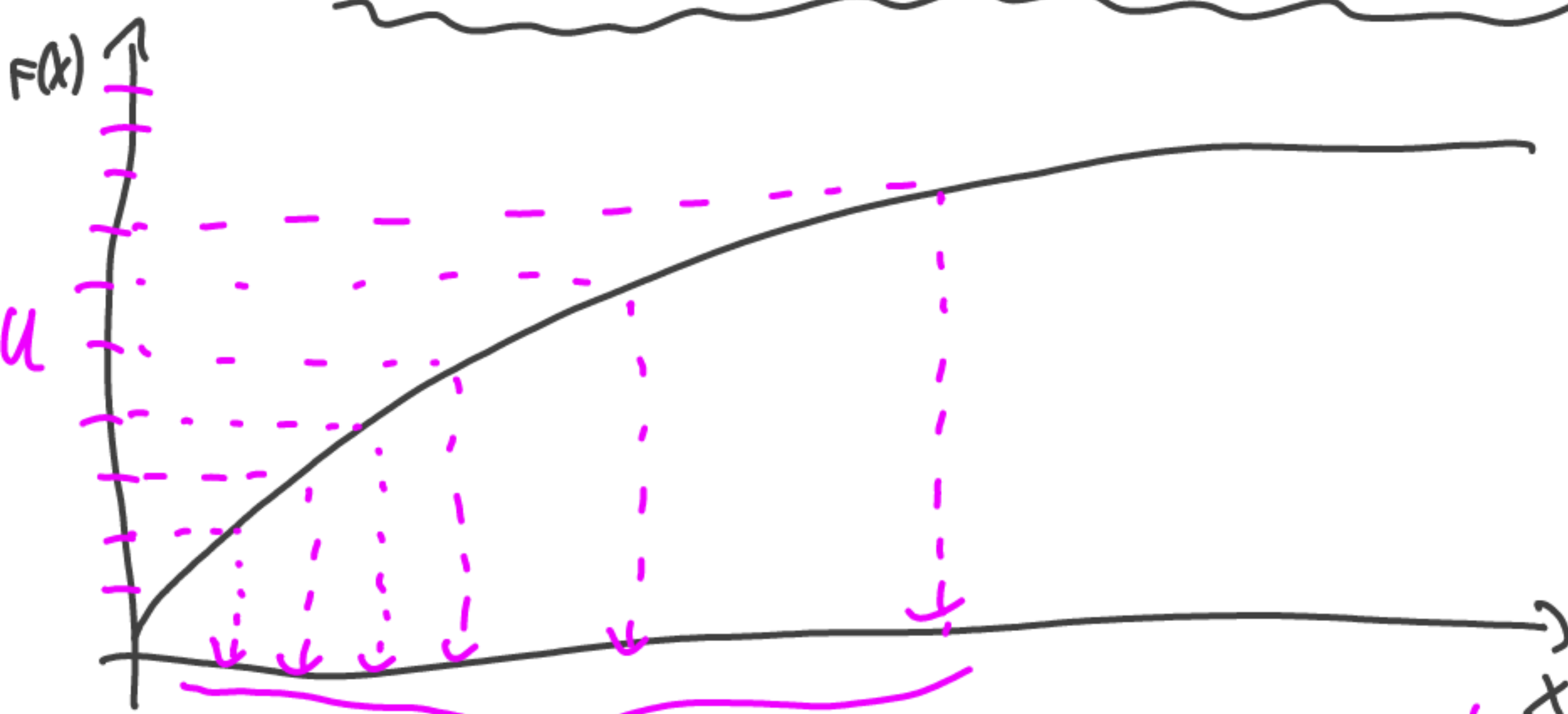
$$\Pr[F^{-1}(U) \leq x] = \Pr[U \leq F(x)] = F(x)$$

\uparrow
F is non-decreasing



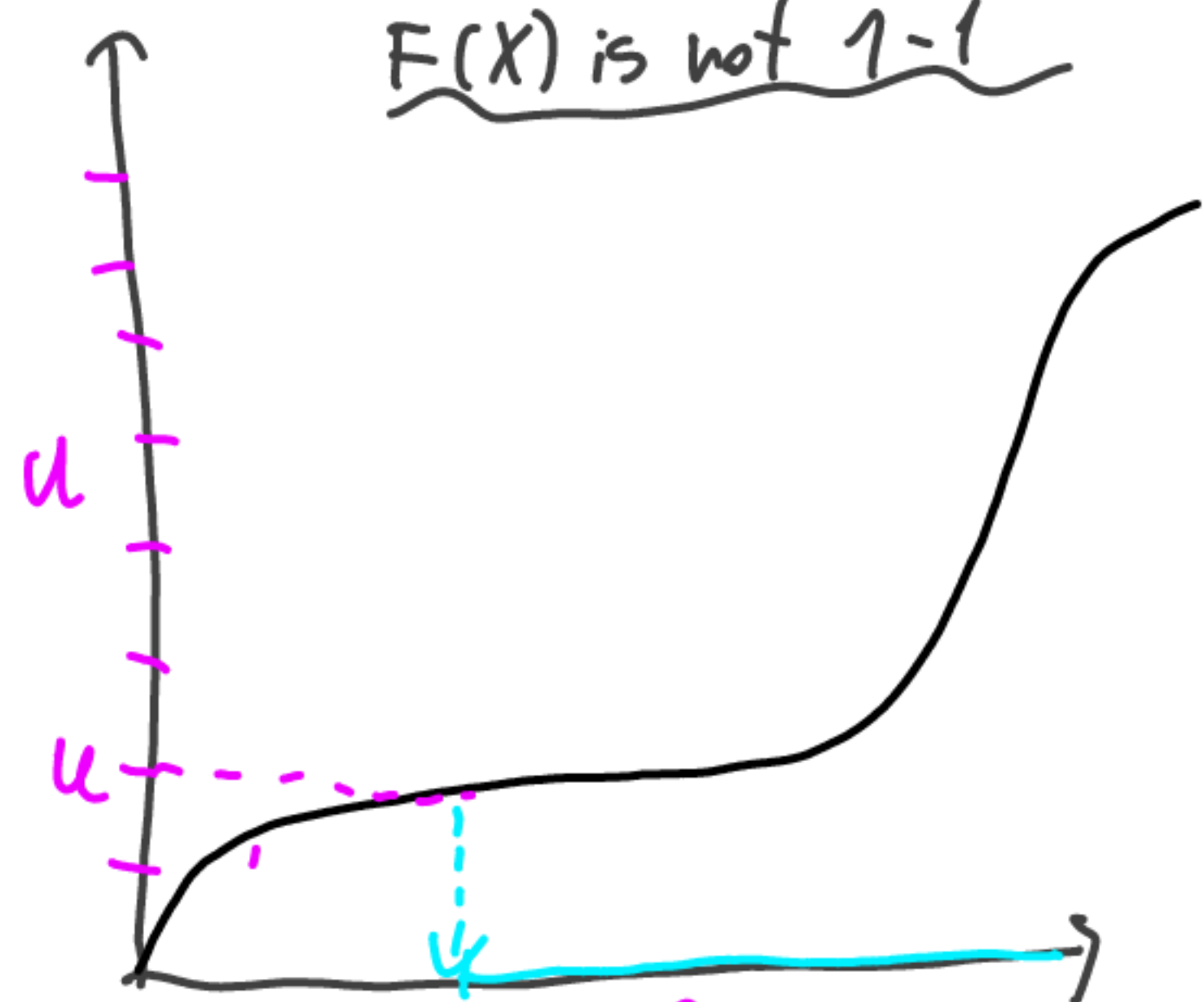
we can use this to obtain any distribution by obtaining elements from Uniform distribution

Examples (why we need infimum in definition of $F^{-1}(x)$)
we assume that $F(x)$ is 1-1 and continuous



there is no longer a dist. it's $F(x)$ dist
 $F^{-1}(u) \sim F(x)$

$F(x)$ is not 1-1

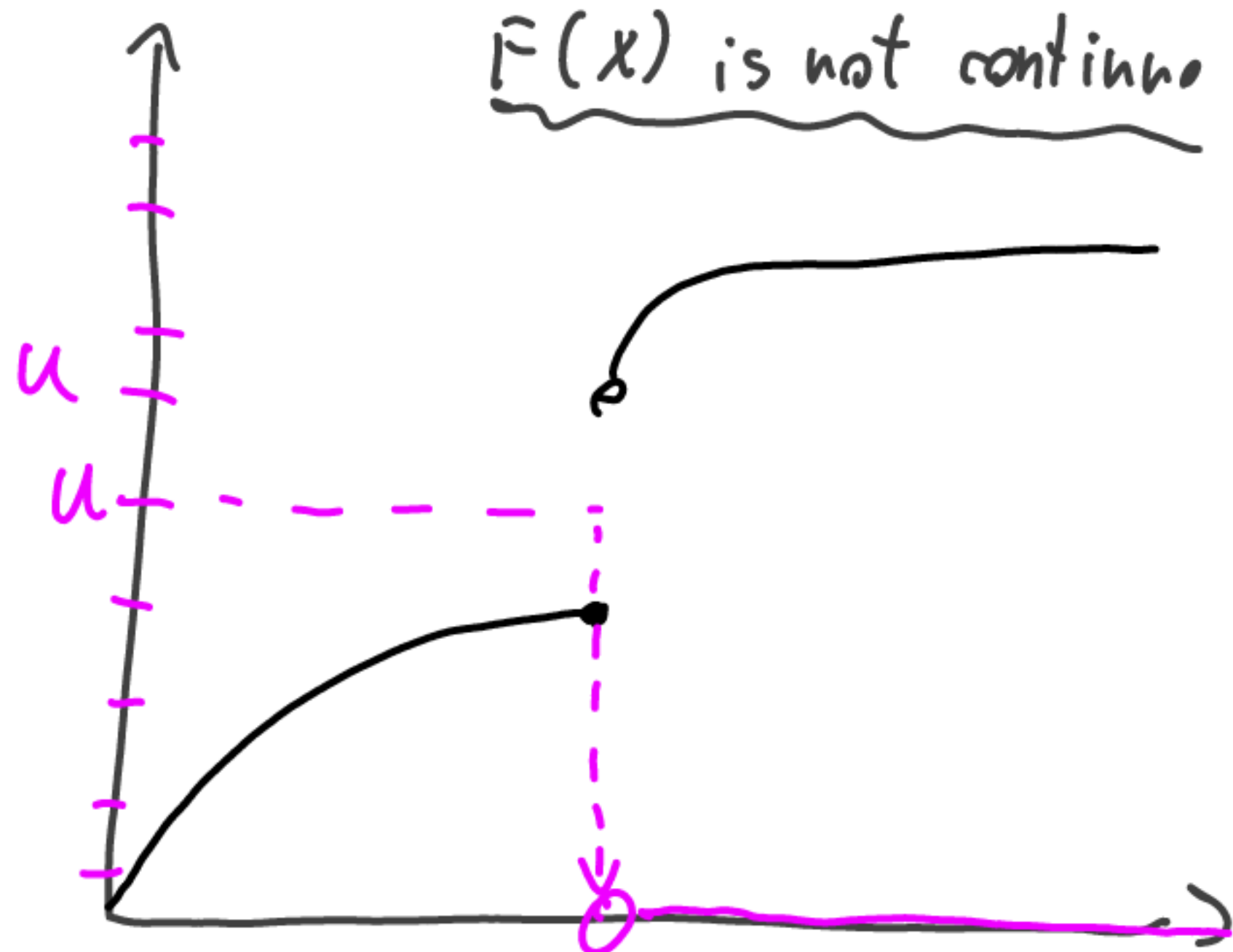


we can choose any point from here.

But from definition we choose x minimum

$$\{x : F(x) \geq u\}$$

$F(x)$ is not continuous



x infimum for $\{x : F(x) \geq u\}$

Generating values from the exponential distribution

$$S \sim \text{Exp}(\lambda), E[S] = \frac{1}{\lambda}, f(x) = \lambda e^{-\lambda x}$$

$$\text{cdf } F(x) = 1 - e^{-\lambda x}, x \geq 0, \lambda > 0$$

we look for inverse of F .

$$u = 1 - e^{-\lambda x} \rightarrow x = \frac{\ln(1-u)}{-\lambda} \rightarrow F^{-1}(u) = \frac{\ln(1-u)}{-\lambda}$$

$$\text{Lemma 1: } U \sim U(0,1) \rightarrow F^{-1}(U) \sim \frac{\ln(1-U)}{-\lambda} \sim \frac{\ln(U)}{-\lambda} \sim \text{Exp}(\lambda)$$

Just in change of sign

The next step is to derive the estimator

Estimator definition



we think about $M[i]$ as results of indep. experiments

$$M_1 = \min\{S_1^{(1)}, S_2^{(1)}, \dots, S_n^{(1)}\}$$

$$M_m = \min\{S_1^{(m)}, S_2^{(m)}, \dots, S_n^{(m)}\}$$

$$M_k = \min\{S_1^{(k)}, S_2^{(k)}, \dots, S_n^{(k)}\}, S_i^{(k)} \sim \text{Exp}(\lambda_i)$$

Lemma 2:

$$\lambda_1 + \lambda_2 + \dots + \lambda_n = \lambda \rightarrow M_k \sim \text{Exp}(\lambda)$$

proof

$$F_{M_k}(x) = P_r[\min\{S_1, \dots, S_n\} \leq x] = 1 - P_r[S_1 > x \wedge S_2 > x \wedge \dots \wedge S_n > x] =$$

$$= 1 - \prod_{i=1}^n P_r[S_i > x] = 1 - \prod_{i=1}^n e^{-\lambda_i x} = 1 - e^{-x \sum_{i=1}^n \lambda_i} = 1 - e^{-x \lambda}$$

" $1 - F(x)$

so for each element of sketch we know what's the distribution.

Conclusion

$$M_k \sim \text{Exp}(\lambda) \rightarrow E[M_k] = \frac{1}{\lambda} \xrightarrow{M_k \approx E[M_k]} \lambda \approx \frac{1}{M_k}$$

from Hyperloglog we learn that harmonic mean has better properties

$$\bar{\lambda} = \frac{m}{\left(\frac{1}{M_1}\right) + \dots + \left(\frac{1}{M_m}\right)}$$

← this estimator with the lowest variance. $= \frac{m}{\sum_{k=1}^m M_k}$

Lemma 3

$$\text{Let } G = \sum_{k=1}^m M_k, \quad M_k \sim \text{Exp}(\lambda)$$

then $G \sim \text{Gamma}(m, \lambda)$, $m \in \mathbb{N}_+$, $\lambda \in \mathbb{R}_+$, $x \geq 0$

Proof (ex 25)

$$f_G(x) = \lambda e^{-\lambda x} \frac{(\lambda x)^{m-1}}{\Gamma(m)}$$

$$, \quad \Gamma(m) = (m-1)! \quad \text{for } m \in \mathbb{N}$$

• it is ok for $m=1$: $f_G(x) = \lambda e^{-\lambda x} \sim \text{Exp}(\lambda)$

• use induction to show that $G_m = G_{m-1} + M_1$ has distribution $\text{Gamma}(m, \lambda)$

tip: take integral

\uparrow
 $\text{Gamma}(m-1, \lambda)$

Lemma 4 $m \geq 2$, $\tilde{\lambda} = \frac{m}{\sum_{k=1}^m M_k}$

we show that $E[\tilde{\lambda}] = \frac{m}{m-1} \lambda$
proof (ex. 26)

$$G = \sum_{k=1}^m M_k \sim \text{Gamma}(m, \lambda) \rightarrow$$

$$\rightarrow E[\tilde{\lambda}] = E\left[\frac{m}{G}\right] \stackrel{(*)}{=} \int_0^{\infty} \frac{m}{x} f_G(x) dx = \dots = \frac{m}{m-1} \lambda$$

$$* E[g(x)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$$

so if we replace m with $m-1$ we get unbiased estimator

Conclusion

lemma 4

$$\tilde{\lambda} := \frac{m-1}{\sum_{k=1}^m M_k}$$

$$\rightarrow E[\tilde{\lambda}] = \lambda$$

best estimator

Lemma 5 $m \geq 3$, $\tilde{\lambda} := \frac{m-1}{\sum_{k=1}^m M_k} \rightarrow SE[\tilde{\lambda}] = \frac{1}{\sqrt{m-2}}$

(ex 27)

proof

$$SE[\tilde{\lambda}] = \sqrt{\text{Var}\left[\frac{\tilde{\lambda}}{\lambda}\right]}$$

$$\cdot \text{Var}[\lambda] = E[\lambda^2] - E[\lambda]^2$$

$$\cdot E\left[\left(\frac{m-1}{G}\right)^2\right] = \dots$$

what we see as

$$\frac{\tilde{\lambda}}{\lambda}$$

