

# Lecture 13: The Data Stream Model

✓ Recap: dimension reduction

- data stream model

- computing 2nd moments

- distinct elements?

← Approx matrix multiplication

→ Approximation  
→ Randomization

$O(\log \log n)$  bits

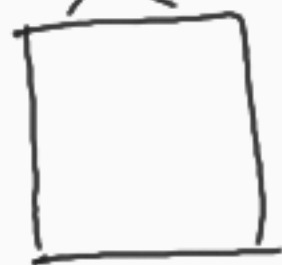
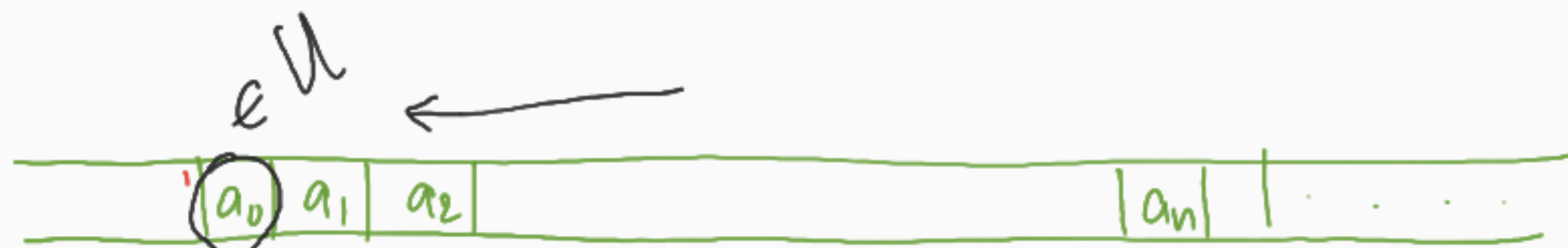
- how many elements?

- distinct?

- sum of elts?

- most frequent?

- median, mean,



little memory

logarithmic

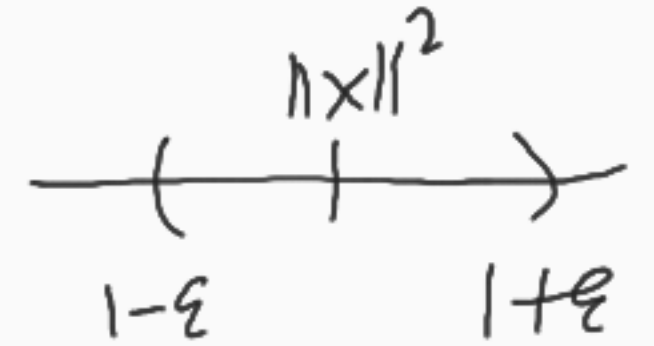
$O(\log \log n)$   
 $\epsilon^{O(n)}$

Recap: JL Lemma:

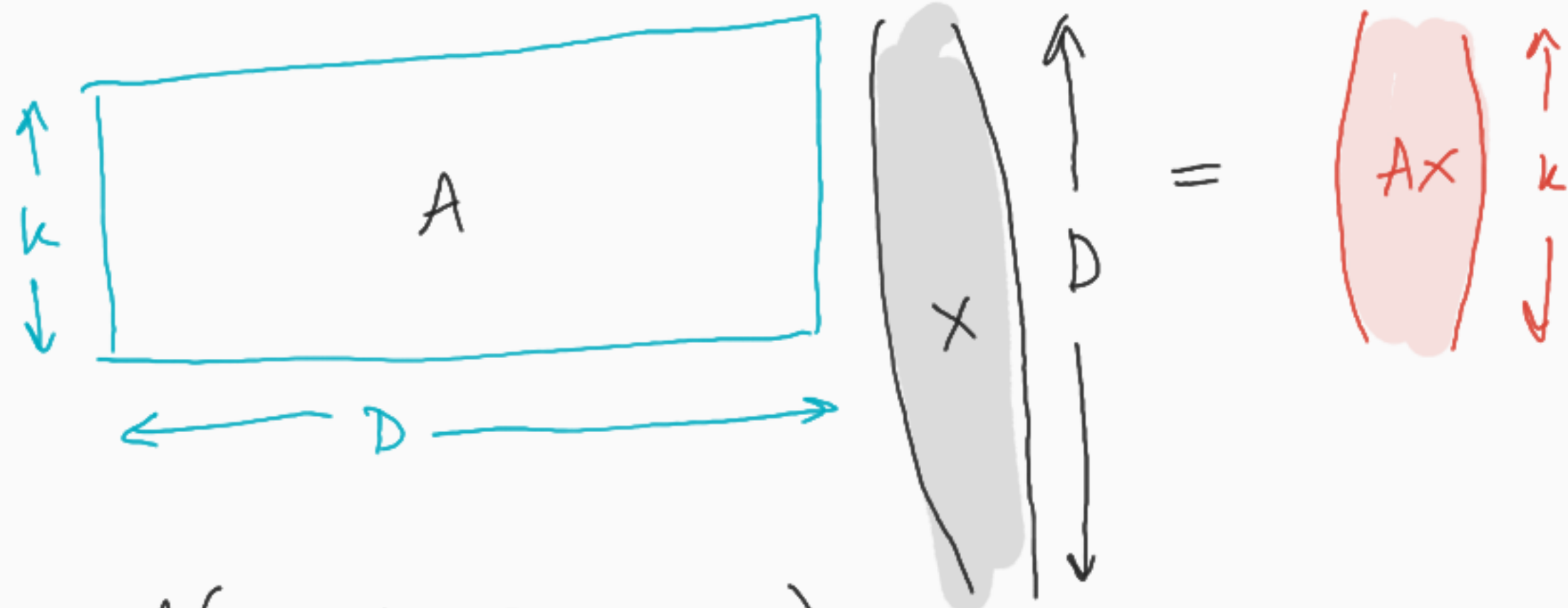
$\rightarrow k \times D$  matrix

gen random matrix  $A = \frac{1}{\sqrt{k}} M$  st  $\forall x \in \mathbb{R}^D$

$$\Pr \left[ \|Ax\|^2 \notin \underbrace{(1 \pm \epsilon)\|x\|^2} \right] \leq \underbrace{\frac{1}{n^2}}_{\delta} \delta$$

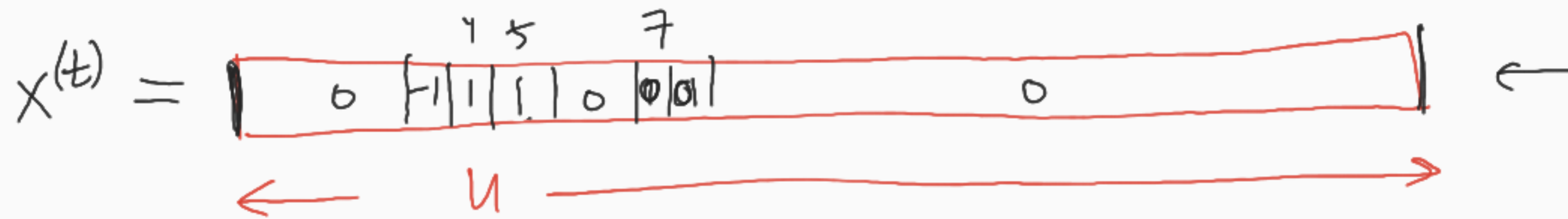
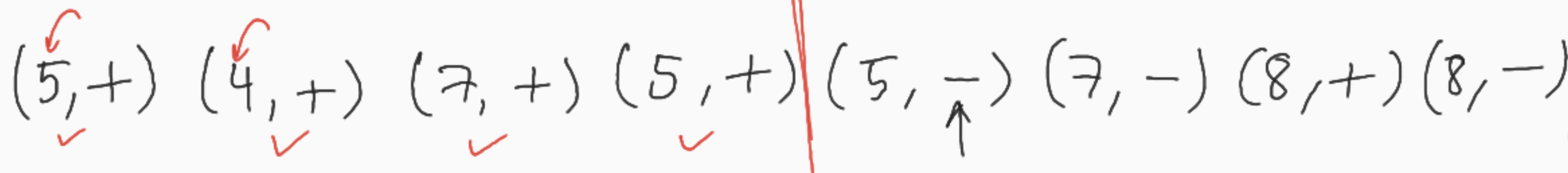


where  $k = c \frac{\log \frac{1}{\delta}}{\epsilon^2} \log \frac{1}{\delta}$



$$A(\underbrace{x_1}_{=} + \underbrace{x_2}_{=} - \underbrace{3x_3}_{=} + \underbrace{5x_4}_{=}) = \underbrace{Ax_1}_{=} + \underbrace{Ax_2}_{=} - \underbrace{3Ax_3}_{=} + \underbrace{5Ax_4}_{=}$$

Computing Second moments ( $F_2$ )



# distinct elts =  $\# \text{nnz}(X) = \|X\|_0$

# active cols =  $\|X\|_1$

$F_2 = \|X\|_2^2$

heavy hitters =  $\|X\|_\infty = \arg \max_i \{X_i\}$

at any time  $t$ , if  $i$  query, should get correct answer w.p

$\geq 1 - \delta$

Computing  $F_2$  [Alon Matias Szegedy]  $C \leftarrow 0$

Fix  $h: U \rightarrow \{-1, 1\}$

4 wise independent hash fn.

f (i, +) then  $C \leftarrow C + h(i)$

(i, -)  $C \leftarrow C - h(i)$

if query then return  $C^2$

---

K wise indep hash family  
of  $U \rightarrow [M]$

$H = \{h: U \rightarrow [M]\}$  is family of fns.

$$\Pr_{h \leftarrow H} [h(a_1) = \alpha_1 \wedge h(a_2) = \alpha_2 \dots \wedge h(a_k) = \alpha_k] = \frac{1}{M^k}$$

Easy to sample from:

$k \log |U|$   
bits of randomness

$\forall a_1 \neq a_2 \neq \dots \neq a_k, \alpha_1, \alpha_2, \dots, \alpha_k \in [M]$   
 $\in U$

$$h(x) = \overbrace{C_0} + \overbrace{C_1}x + \overbrace{C_2}x^2 + \overbrace{C_3}x^3$$

↑

↑

↑

↑

4-wise indep

all arithmetic over finite field

choose uniformly from  $\mathbb{F}_{2^n}$

$U \rightarrow \cancel{2^m = M}$

$\parallel$   
 $2^u \rightarrow 2^u$

↓  
 drop highest  $u-m$  bits

$\Rightarrow$  4-wise indep hash fn  
 from  $U \rightarrow [M]$

total randomness =  $4 \lg |U|$

Lemma 1:  $E[C^2] = F_2$

$$E\left[\left(\sum_i x_i h_i\right)^2\right]$$

$$= E\left[\left(\sum_i x_i h_i\right)\left(\sum_j x_j h_j\right)\right]$$

$$= E\left[\sum_i x_i^2 h_i^2 + \sum_{i \neq j} x_i x_j h_i h_j\right]$$

$$= \sum_i x_i^2 + \sum_{i \neq j} x_i x_j \underbrace{E[h_i]}_0 \underbrace{E[h_j]}_0$$

$$= \sum_i x_i^2 = F_2$$

$C \leftarrow \sum_i x_i h_i$

$$\text{Var}(C^2) = E[(C^2)^2] - \underbrace{E[C^2]^2}$$

$E[C^4]$

$$E\left[\left(\sum_i x_i h_i\right)^4\right]$$

$$E\left[\sum_i x_i^4 h_i^4 + \sum_{i \neq j} x_i^2 x_j^2 h_i^2 h_j^2\right]$$

$$+ \sum_{i \neq j \neq k} x_i^2 x_j x_k h_i^2 h_j h_k$$

$$+ \sum_i x_i^3 x_j h_i^3 h_j$$

$$4 \sum_{i \neq j} x_i^2 x_j^2 = \sum_i x_i^4 + c \sum_{i \neq j} x_i^2 x_j^2 - (\sum_i x_i^4 + 2 \sum_{i \neq j} x_i^2 x_j^2) =$$

$$\sum_i x_i^4 + c \sum_{i \neq j} x_i^2 x_j^2 - (\sum_i x_i^2)^2$$

Lemma 1:  $E[C^2] = F_2 = \sum_i x_i^2$

Lemma 2:  $\text{Var}(C^2) = 4 \sum_{i \neq j} x_i^2 x_j^2$

$$Pr[|C^2 - \mu| \geq \epsilon \mu] \leq \frac{\text{Var}(C^2)}{\epsilon^2 \mu^2}$$

Padnuti  
Chebyshev's Ineq.

$C_1, C_2, \dots, C_k$

means  $\leftarrow \frac{\sum_{i=1}^k C_i^2}{k}$

mean  $F_2$   
var  $\leq \frac{4}{k}$

$X_1, X_2, \dots, X_k$  iid  
 $(\mu, \sigma^2)$

$$= \frac{4 \sum_{i \neq j} x_i^2 x_j^2}{k \epsilon^2 \left( \sum_i x_i^2 \right)^2}$$

$$= \left( \frac{\sum_i x_i^2}{k} \right)$$

$$\leq \frac{4}{\epsilon^2} k \text{ (bigger than 1!)} \quad \text{☹}$$

$$= \delta \Rightarrow k = \frac{4}{\epsilon^2 \delta}$$

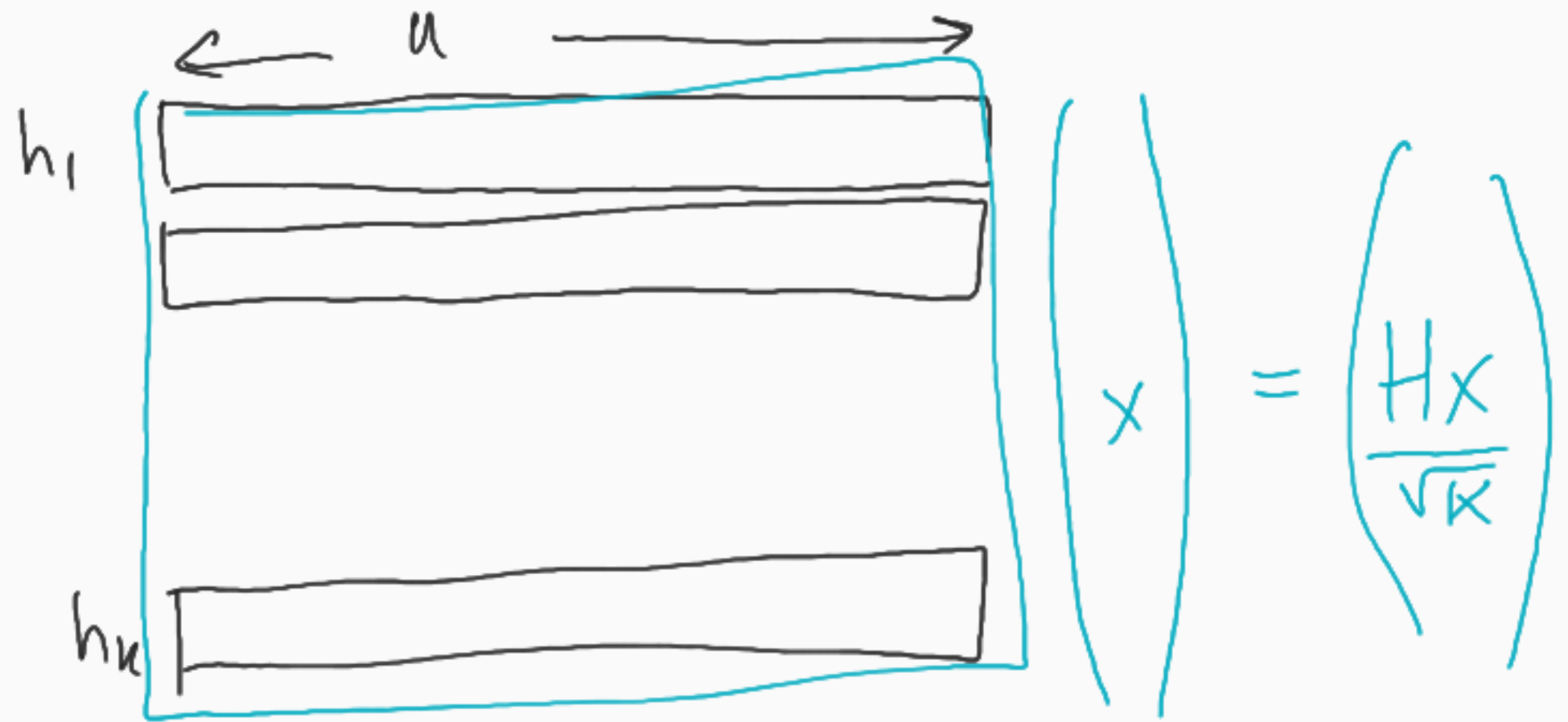
$\mu, \sigma^2/k$

$h_1$   
 $h_2$

$$C_1 = \sum_i x_i h_{1i}(x)$$
  
$$h_2 \cdot x$$

$h_k$

$$C_k = \sum_i x_i h_{ki}(x)$$



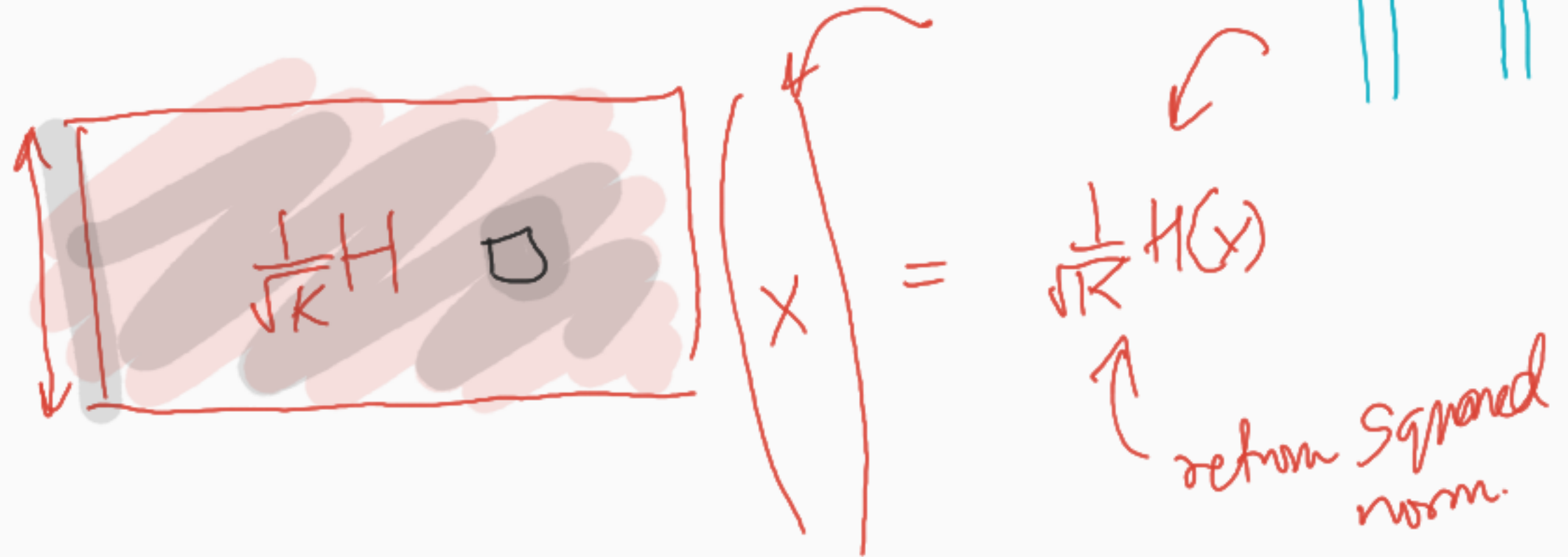
reduced

$$\left( \frac{\sum_{j=1}^k C_j^2}{K} \right)$$

large space

sm. space

$$\frac{1}{\sqrt{K}} H$$



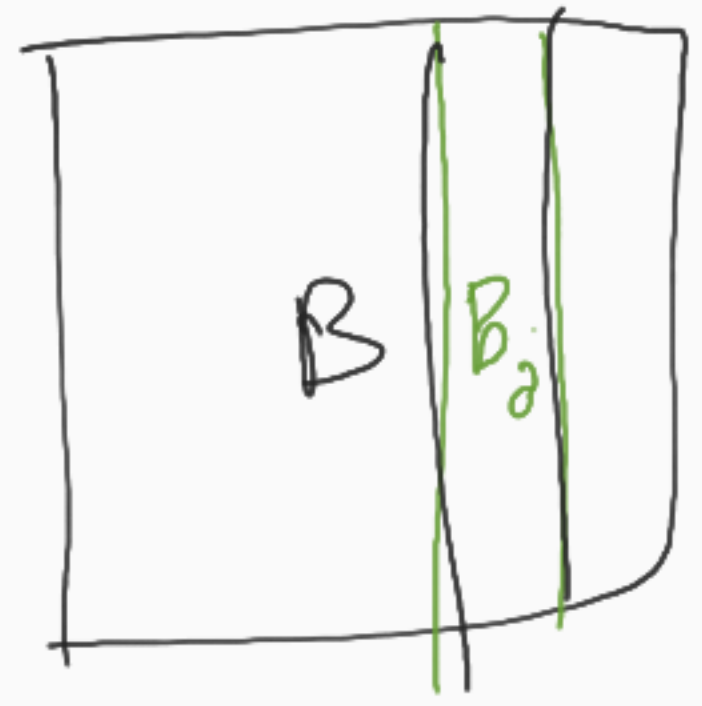
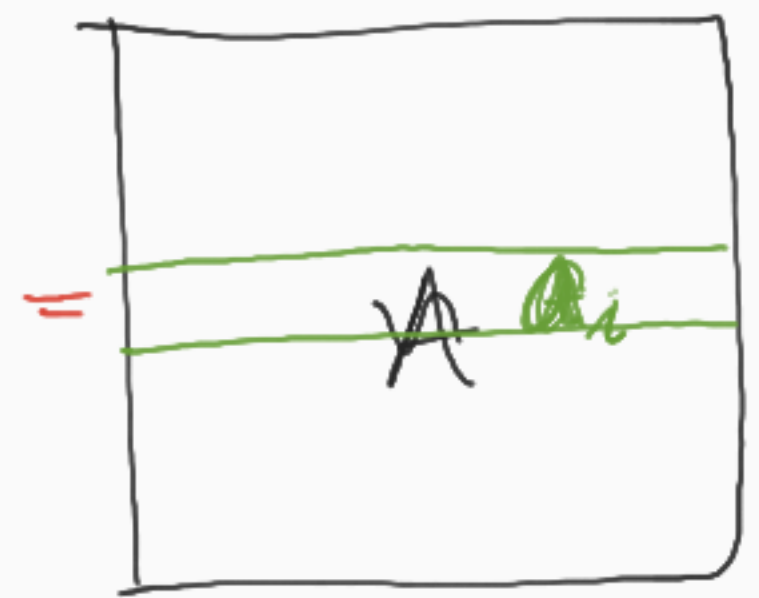
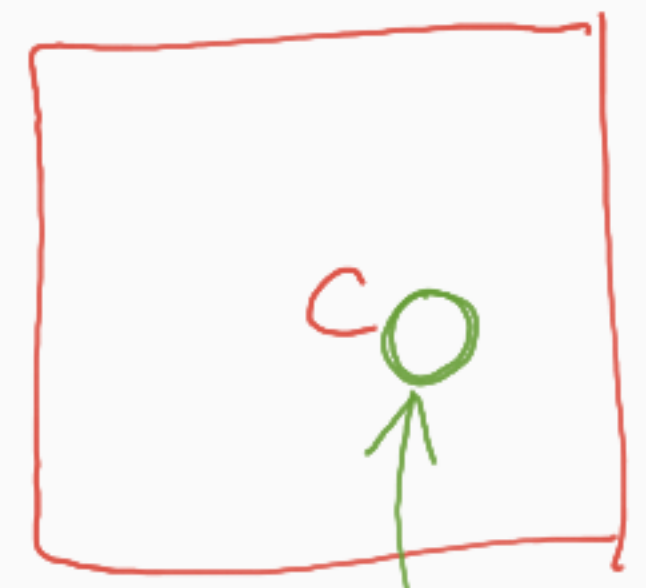
JL:  $K = \frac{105 \cdot 1/8}{\epsilon^2}$

Chelby:  $\frac{1/8}{\epsilon^2}$

$$\| \cdot \|_2^2$$



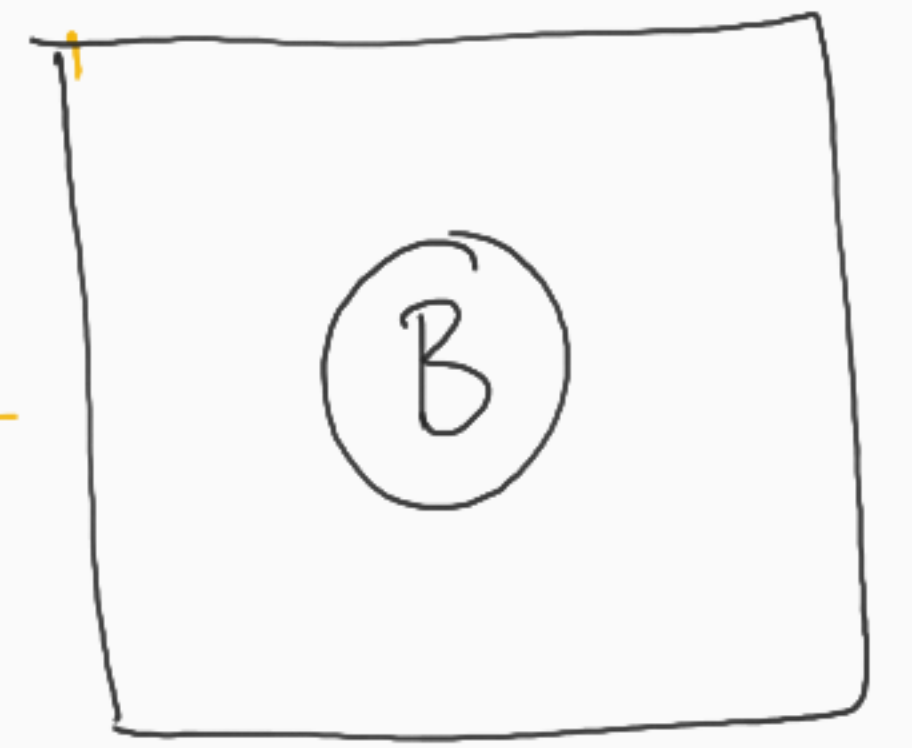
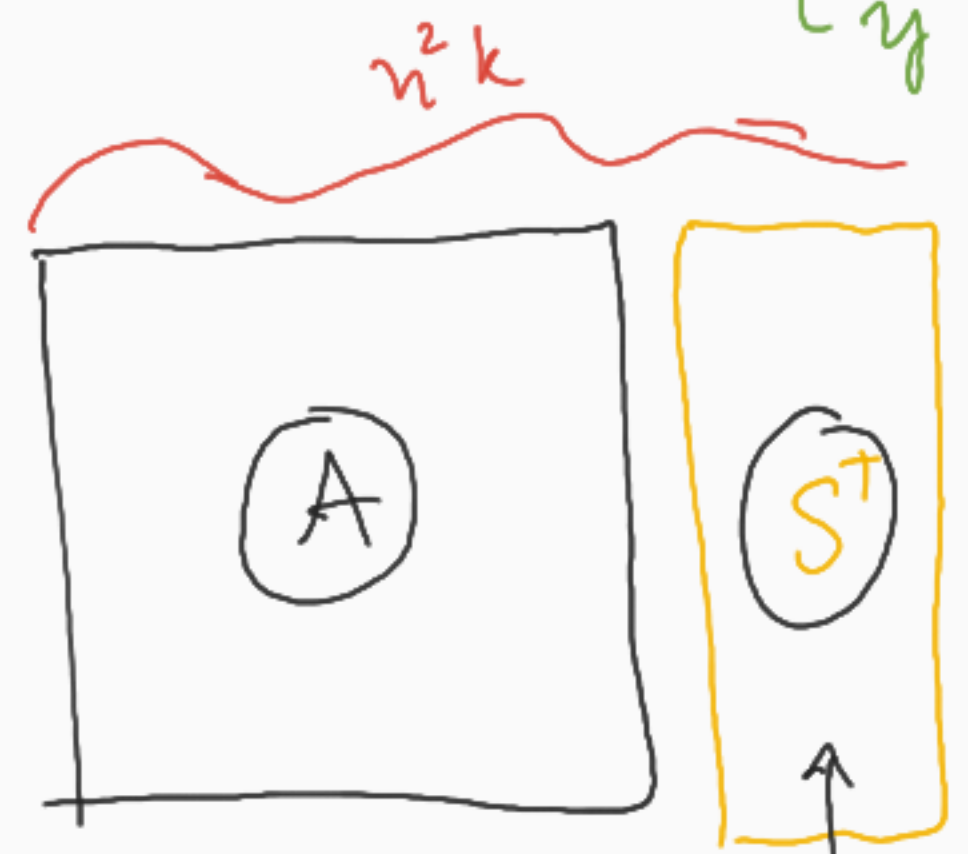
AB



$n^3$

$C_{ij} = a_i \cdot B_j$

approx  
 $C \approx$



$nk^2$

JL matrix  
or Chebyshev matrix

$k \times n$







