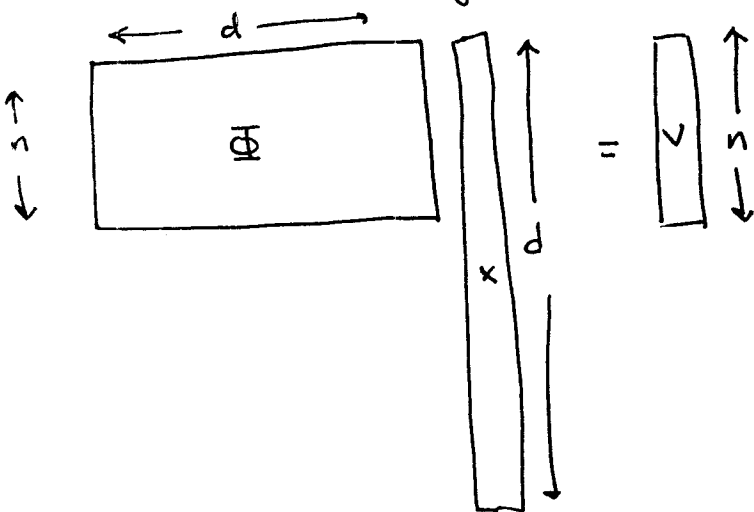


What makes sublinear algs. so fast?

① statistical signal recovery ← WHAT problem



observe signal $x \in \mathbb{R}^d$, linear measurements Φ
 record observations $v \in \mathbb{R}^n$. Retain Φ, v .

Output: info. about x

ℓ_p norm, median (order statistics), x itself,
 repr of x in basis/dictionary, $\{0,1\}$ decision about x ,
 $\langle x, q \rangle$ with arbitrary q

Example of sublinear alg:

Thm: [GSTV '07]

[these types of results have a lot of moving parts — let's try to understand pieces.]

With high prob., a random matrix Φ with $m \text{poly}(\log(d/\epsilon))$ rows ~~is~~ is a good meas. matrix. \exists alg. s.t. given

$\Phi x = v$ and $m = \text{sparsity}$, alg. outputs \hat{x} (m non-zero entries)
 s.t. (\hat{x} has optimal error guarantees.)

$$\|x - \hat{x}\|_2 \leq \|x - \hat{x}\|_1 + \frac{\epsilon}{\sqrt{m}} \|x - \hat{x}\|_2$$

in time $\Theta(m^2 \epsilon^{-4} \text{poly} \log d)$.

② In these lectures, we'll walk through 2 canonical stat. signal recovery algs. SLOWLY!

① Combinatorial Group Testing / Hashing

input: x binary vector, length d
 m non-zeros

output: m pos's of spikes (or reasonable fractⁿ)

② 2nd freq. moment, energy of stream

input: $x \in \mathbb{R}^d$ in a data stream

output: $\|x\|_2^2 = \sum_{i=1}^d |x(i)|^2$ (or close approxⁿ)

③ Models / Applications ← WHERE do these probs. arise?

(i) "Compressed Sensing"

(ii) Data streams (IP networks, phone call records, transactions streams)

items arrive sequentially

define a signal $x: [1, \dots, d] \rightarrow \mathbb{R}$ implicitly

rep. transactions

$\langle i, u(i) \rangle$

$i \in [1, \dots, d]$ index

$u(i) = \text{update to } x(i) \Rightarrow x(i) += u(i).$

domain vals presented in arbitrary order

updates are positive OR negative ~~or both~~

$d = \text{domain size HUGE!}$

Alg'ic challenges in data streams

size, volume enormous + most hardware simple, cheap!

ex: OC192 backbone link 9.9 Gbits/sec

IP pkt hdrs = 30 million items/sec

avg. flow = 10 pkts \rightarrow lasts 3.3×10^{-7} sec.

(A) Computational resources

WHAT are we charged for?
i.e., what could be sublinear?

- (1) SPACE - have to store Φ, v
 need auxiliary, working memory
 # bits, precision

all require RAM, physical space, power

ex: all pkt hdrs on OC192 for 1 month = 1000s PBytes

- (2) TIME - have to output answer
 operate on Φ, v

post-processing: apply Φ

CS / data stream models. {

- (a) NOT charged to collect v
 except 1 entry = unit cost
- (b) charged PER UPDATE to x

ex: router has to process each pkt quickly

(3) COMMUNICATION - send Φ, v to central location

(4) RANDOMNESS -- Φ is a random matrix
generate AND store Φ , + random bits

More practically, hardware = SWEPT metrics.
size, weight, energy, power, time

⑤ Sublinear Algs.

Informally, resource usage grows slower than size of input. If double input size, use LESS than 2x more resources to compute answer.

ex: input size = d

space/time $O(d^\alpha)$ for $\alpha < 1$

$O(\log d)$

polylog(d) = $\log^c(d)$ for some c .

Consequences of sublinear resources (or HUGE data!)

- can't output x itself
- can't take too many meas.
- can't store Φ explicitly
- can't do bit stuffing
- can't use lots of indep. bits
- post process fast
- one pass

HARD :

data stream of active ph. #'s
1 million phone #'s
see 100K records
only 4 active at time
you look
where are they?!

CGT

- Goes back to WWII, 2 economists in OPA wanted an efficient method for testing draftees for syphilis

- (1) draw blood from each person
- (2) pool blood samples into groups of 5
- (3) use a single test on each group

OBSERVE

- (1) If no one in pool has antigen, then test negative
- (2) If ≥ 1 people have antigen, then pool does
and test is positive

↳ can then test each individually

- \therefore If no one has antigen, save 4 tests (over all 5)
- If ≥ 1 —————, waste 1 test (over all 5)

Def. (CGT) Have a universe of d items, m are "defective"

GOAL: Construct a collection of tests (a design) to minimize the # tests needed to find defective set for worst case input.

Not specified: (1) type of tests

binary / not

linear / not

adaptive / not

det. / prob.

(2) algorithm

WANT: linear, binary, non-adaptive tests

good & fast algs.

prob. ok.

WARM UP
EXAMPLE:

$x =$ binary signal, length d , 1 sparse

Defⁿ: Let B_1 be a binary matrix with $n = \lg d$ rows and d cols. The i^{th} col. of B_1 is i written in

binary. BT tester matrix

$\bar{B}_1 =$ bit tester with add'l top row $= 1$.

$$\begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$$

Algorithm: ① Compute $B_1 x = v$

② output $v = \text{pos}^n$ of defect (in binary)

Analysis: ① space for $v = \lg d$

② alg. is trivial, takes time $\lg d$ to output

③ don't have to store B_1 explicitly

receive update in $\text{pos}^n i$,

generate i in binary,

for each entry in i , multiply update

accumulate meas.

PROBLEM #1: let x be m -sparse binary vector length d

GOAL: pos^n s of m non-zero entries of x .

lemma: There is a ^{random} measurement matrix Φ and an ~~algorithm~~ alg

that, for each x of length d with m defects, returns a

list of at most $10m$ items that contains at least $\frac{m}{10}$

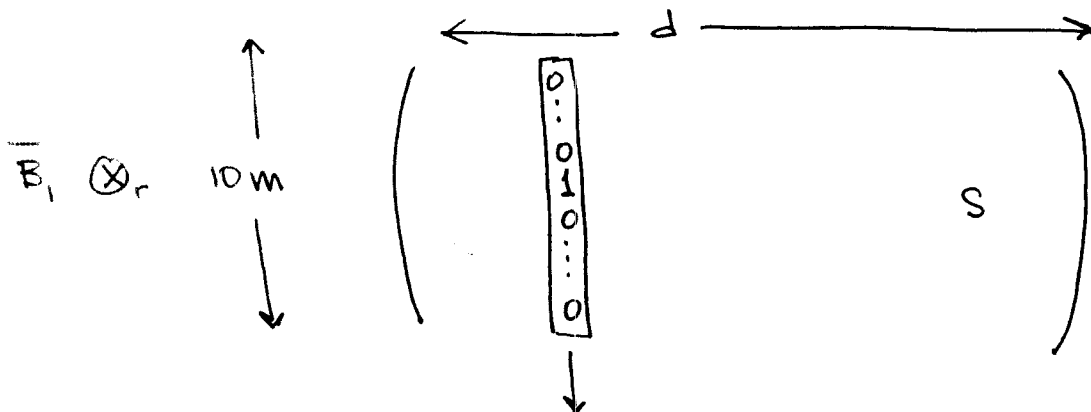
defects with prob. $\geq \frac{8}{9}$. The alg. takes time $O(m \lg d)$

[and uses $O(m \lg d)$ measurements.]

Defⁿ: $B \otimes_r S = M =$ row tensor product of 2 matrices.

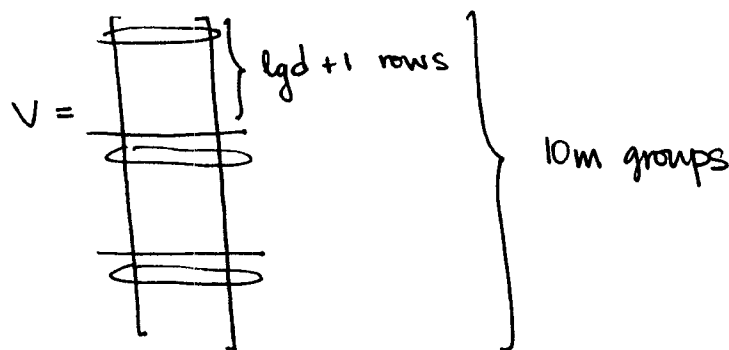
Each row of $M =$ elt. wise prod. of 1 row in B and 1 row in S .
 [loop over all rows of B, S to get elt. wise prod.
 over all pairs of rows from B, S .]

Measurement Φ : $\Phi = \bar{B}_1 \otimes_r S =$ binary matrix with
 $10m$ lgd rows, d cols.



1 non-zero entry per col.
 posⁿ chosen indep^{ly}, unif^{ly} at random
 i.e., choose row with prob = $\frac{1}{10m}$

Algorithm:



List = \emptyset

for each of $10m$ blocks

check 1st entry, if non zero

then lgd rows give posⁿ of spike

add to list.

output list.

Analysis:

- observe that each row of S selects entries of x at random
there are really only m entries we care about

$\Rightarrow S$ is placing m balls into $10m$ blets randomly

\therefore need to analyze balls/bins process

[enemy = collisions = balls in same bin]

- Assuming that a row of S isolates a spike,
then \bar{B}_i finds it with the right label.

Note top row just sums all entries in bkt.

- Analyze balls/bins and collisions

Consider a particular spike x_j

it's in some group k ; $S(x_j) = k$ by abuse of notation

- Conditional prob. that another spike x_ℓ is ~~ALSO~~
assigned to k by S

$$\Pr(S(x_\ell) = k \mid S(x_j) = k) = \frac{\Pr(S(x_\ell) = k \cap S(x_j) = k)}{\Pr(S(x_j) = k)}$$

\swarrow independence of balls/bins assign.

$$= \Pr(S(x_\ell) = k) = \frac{1}{10m}$$

- Let $Y_\ell = X_k(S(x_\ell)) = \begin{cases} 1, & \text{prob } \frac{1}{10m} & [\text{if } x_\ell \text{ lands in } k] \\ 0, & \text{prob. } 1 - \frac{1}{10m} & [\text{o.w.}] \end{cases}$

Then let $Y = \#$ other spikes in group k

$$= \sum_{l=1}^{m-1} Y_l$$

and $\mathbb{E}(Y) =$ expected $\#$ other spikes in k

$$= \mathbb{E}\left(\sum_{l=1}^{m-1} Y_l\right) = \sum_{l=1}^{m-1} \mathbb{E}(Y_l)$$

linearity of expectation.

$$= \frac{m-1}{10m} = \frac{1}{10} - \frac{1}{10m} < \frac{1}{10}$$

i.e., on avg. $\#$ other spikes assigned is $\frac{1}{10}$

$$\Rightarrow \Pr(Y \geq 10 \cdot \mathbb{E}(Y)) = \Pr(Y \geq 1)$$

$$\leq \frac{\mathbb{E}(Y)}{10 \cdot \mathbb{E}(Y)} = \frac{1}{10}$$

$$\therefore \Pr((Y \geq 1)^c) = \Pr(\text{no other spike in } k) \geq \frac{9}{10}$$

\Rightarrow each defect is isolated in its own group

with prob. $\geq \frac{9}{10}$ (and it fails to be isolated prob $\leq \frac{1}{10}$)

• Let's look at all the balls now, not just one.

$$\text{Let } Z_j = \begin{cases} 1, & \text{if } x_j \text{ NOT isolated, prob } \frac{1}{10} \\ 0, & \text{if } x_j \text{ isolated, prob } \frac{9}{10} \end{cases}$$

$$\text{and set } Z = \sum_{j=1}^m Z_j = \# \text{ spikes NOT isolated}$$

$$\mathbb{E}(Z) = \sum_{j=1}^m \mathbb{E}(Z_j) \leq \frac{m}{10} = \text{expected \# of spikes NOT isolated}$$

$$\text{again, } \Pr\left(Z \geq \frac{9m}{10}\right) \leq \frac{\frac{m}{10}}{\frac{9m}{10}} = \frac{1}{9}$$

i.e., # of balls we fail to isolate exceeds $\frac{9m}{10}$

with prob. $\frac{1}{9}$ and

$$\Pr\left(\left(Z \geq \frac{9m}{10}\right)^c\right) \geq \frac{8}{9}$$

with prob. $\frac{8}{9}$ we isolate at least $\frac{m}{10}$ defects.

Note that total # groups = $10m$, so the # false positives is under control. Actually, in this simple example, we

KNOW when we have a false positive BUT... you can

generalize this problem and still use same framework. //

Extensions

- ① iterate to get more (all) spikes
- ② increase success prob. with repeated trials, high enough prob. to get ALL m -sparse sig.

STILL NEED

- ① small-space
- ② limited randomness

Pairwise indep. and hash functions

Def'n: A set of rvars X_1, X_2, \dots, X_n is k -wise indep.

if, for any subset $I \subseteq [1, \dots, n]$ with $|I| \leq k$ and

for any vals. $x_i, i \in I$,

$$\Pr \left(\bigcap_{i \in I} X_i = x_i \right) = \prod_{i \in I} \Pr \left(X_i = x_i \right).$$

Example: pairwise indep. mod prime p

Construct Y_0, Y_1, \dots, Y_{p-1} pairwise indep. rvars that are uniform over $\mathbb{Z}_p = \{0, 1, \dots, p-1\}$.

Use 2 indep. uniform vals. $X_1, X_2 \in \mathbb{Z}_p$.

$$\text{Let } Y_i = (X_1 + i X_2) \bmod p.$$

Lemma: Y_0, \dots, Y_{p-1} are 2-wise indep. unif. vars over \mathbb{Z}_p .

proof: (1) Each Y_i unif'ly distributed over \mathbb{Z}_p .

Given X_2 , possibilities for $X_1 \Rightarrow p$ distinct vals.
for $Y_i \pmod p$.
unif. too. ✓

(2) Consider Y_i, Y_j and $a, b \in \mathbb{Z}_p$. Want

$$\Pr(Y_i = a \cap Y_j = b) = \frac{1}{p^2}$$

$$\left. \begin{matrix} Y_i = a \\ Y_j = b \end{matrix} \right\} \iff \begin{cases} X_1 + iX_2 = a \pmod p \\ X_1 + jX_2 = b \pmod p \end{cases}$$

↑
syst. of eqns
2 unknowns, 2 eqns
 $\therefore 1$ soln!

$$X_2 = \frac{(b-a)}{(j-i)} \pmod p \quad X_1 = a - \frac{i(b-a)}{(j-i)} \pmod p$$

X_1, X_2 are unif iid over $\mathbb{Z}_p \therefore$ done. ✓

Hash functions

We've defined Y_0, \dots, Y_{p-1} as a collⁿ of vars but we can also view them as a function $Y: \mathbb{Z}_p \rightarrow \mathbb{Z}_p$ which takes elts. or indices in \mathbb{Z}_p and maps them randomly (pwise indep'ly) to other elts. in \mathbb{Z}_p .

γ is called a hash function - it takes a collection of items $\{0, 1, \dots, p-1\}$ and hashes them into \mathbb{Z}/p .

$$\gamma(0) = \gamma_0$$

$$\gamma(1) = \gamma_1$$

\vdots

$$\gamma(p-1) = \gamma_{p-1}.$$

Defⁿ: Let U be a universe with $|U| \geq \frac{n}{k}$ and let

$V = \{0, 1, \dots, \frac{n-1}{k}\}$. A family of hash functions

$\mathcal{H}: U \rightarrow V$ is k -universal if for any elts.

x_1, x_2, \dots, x_k and for a hash function h chosen unif^{ly} at random from \mathcal{H} , we have

$$\Pr(h(x_1) = h(x_2) = \dots = h(x_k)) \leq \frac{1}{n^{\frac{k-1}{k}}}.$$

We want $|U| = d \gg n$; i.e., hash ~~into~~ a large # items into a smaller # bins.

and, frequently, just 2-universal hash functions, as we care about collisions.

Example: 2-universal hash function

Let $U = \{0, 1, \dots, d-1\}$ universe

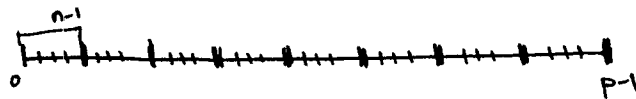
$V = \{0, 1, \dots, n-1\}$ $n < d$ (= range of hash)

choose prime $p \geq d$.

consider family $\mathcal{H}: U \rightarrow V$ given by

$$\mathcal{H} = \{h_{a,b} \mid a \in \mathbb{Z}_p^*, b \in \mathbb{Z}_p\} \quad a \neq 0!$$

$$h_{a,b}(x) = ((ax + b) \bmod p) \bmod n$$



Lemma: \mathcal{H} is 2-universal.

proof: Need to look at how many of the functions $h_{a,b}$ have collisions; i.e., as a function of a, b for $x_1 \neq x_2$ pairs x_1, x_2 , how many have

$$h_{a,b}(x_1) = h_{a,b}(x_2) ?$$

omit ...

To apply this to our CGT problem —

recall that $n = 10m$ (hash into $10m$ buckets)

and that we needed

$$\Pr(S(x_\ell) = k \mid S(x_j) = k) = \frac{1}{10m} = \frac{1}{n}$$

- If S is one of our hash functions, then we're done.
- \therefore pick $a \in \mathbb{Z}_p^*$, $b \in \mathbb{Z}_p$ unif'ly indep'ly at random
 $\hookrightarrow \mathcal{O}(\lg p) = \mathcal{O}(\lg d)$ bits of randomness

- to compute

$$S(x_\ell) = h_{a,b}(x_\ell) = [(ax_\ell + b) \bmod p] \bmod n.$$

ℓ_2 norm of data stream

Alon, Matias, Szegedy '96.

input: data stream $\langle i, u(i) \rangle$ defines $x \in \mathbb{R}^d$ implicitly

GOAL: compute $\sum_{i=1}^d |x(i)|^2 = \|x\|_2^2$ highly eff'ly.

Thm: (AMS '96)

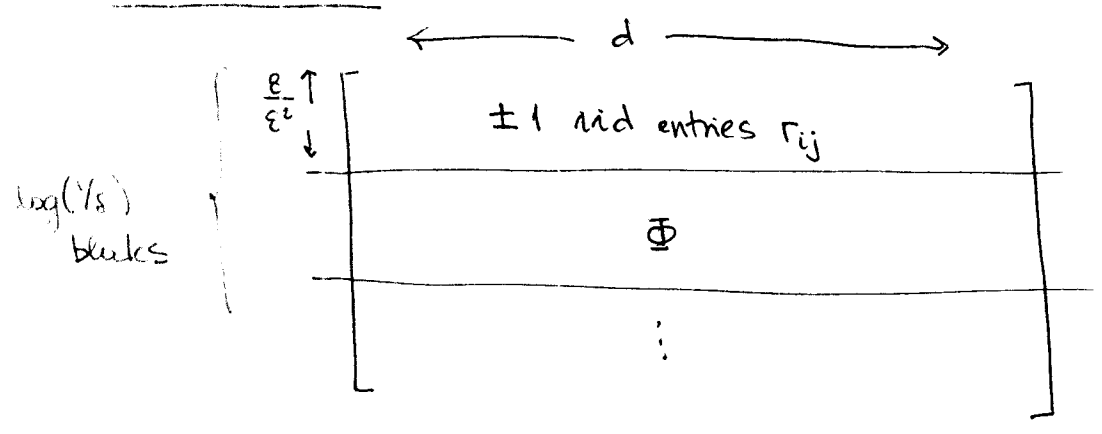
There is a streaming alg. which takes as input a stream of records (defining x), and for each x , with prob. $1 - \delta$, produces an estimate Z of $\|x\|_2^2$ which sat.

$$(1 - \epsilon) \|x\|_2^2 \leq Z \leq (1 + \epsilon) \|x\|_2^2.$$

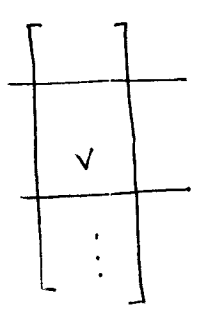
The alg. uses $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$ space and both returns the estimate and process each stream record in time $O\left(\frac{1}{\epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$. The alg. uses $O(\log d)$ random bits.

proof:

Measurements:



Algorithm:



each row j of v is an estimator $X_j = \sum_{i=1}^d r_{ij} x(i)$

Let $Z_j = X_j^2 = \left(\sum_{i=1}^d r_{ij} x(i)\right)^2$

for each block k of $\frac{8}{\epsilon^2}$ estimators Z_j , compute

$Y_k = \frac{\epsilon^2}{8} \sum_j Z_j$ average Z_j on each block.

output $W = \text{Median}_k Y_k = \text{median over all blocks.}$

Analysis:

(1) Z_j is unbiased est.

$$\begin{aligned}\mathbb{E}(Z_j) &= \mathbb{E}(X_j^2) = \sum_i \mathbb{E}(r_{ij}^2) x(i)^2 + \sum_i \sum_{i \neq l} \cancel{\mathbb{E}(r_{ij}) \mathbb{E}(r_{lj})} x(i)x(l) \\ &= \|x\|_2^2\end{aligned}$$

(2) 2nd moment of Z_j

$$\mathbb{E}(Z_j^2) = \mathbb{E}(X_j^4) = \sum x(i)^4 + 6 \sum_{i \neq l} x(i)^2 x(l)^2$$

$$\begin{aligned}\Rightarrow \text{Var}(Z_j) &= \mathbb{E}(Z_j^2) - \mathbb{E}^2(Z_j) \\ &= 4 \sum_{i \neq l} x(i)^2 x(l)^2 \\ &\leq 2 \|x\|_2^4\end{aligned}$$

(3) $Y_k = \text{avg. of } \frac{8}{\varepsilon^2} \text{ unbiased est.} \Rightarrow Y_k \text{ also unbiased}$

$$\mathbb{E}(Y_k) = \|x\|_2^2$$

$$\text{Var}(Y_k) = \frac{\varepsilon^2}{8} \text{Var}(Z_j) \leq \frac{\varepsilon^2 \|x\|_2^4}{4}$$

(4) Chebyshev / Markov

$$\Pr(|Y_k - \mathbb{E}(Y_k)| \geq \varepsilon \|x\|_2^2) \leq \frac{\text{Var}(Y_k)}{\varepsilon^2 \|x\|_2^4} \leq \frac{1}{4}$$

is) $W = \text{median of } \log(1/8) \text{ vars } Y_k$

if $\left|W - \|x\|_2^2\right| > \varepsilon \|x\|_2^2$, then half of Y_k s

are too large, but we expect this to happen
for only $1/4$ of Y_k s.

Chernoff \Rightarrow bad events happen for $> \frac{1}{2}$ Y_k s

with prob. expon'ly low in #vars = $\Theta(\log 1/8)$. \checkmark

Discussion:

(1) Did we need full independence? NO!

just need 4-wise to bound $\mathbb{E}(Z_j^2)$

(2) How do we construct 4-wise indep. ± 1 (or 0/1)
vars quickly and compute eff'ly?

5-wise indep. rvars and error correcting codes

Thm: Suppose $d = 2^k - 1$ and $S = 2t + 1$ (so $t = 2$).

Then there exists a symmetric prob. space Ω of size $2(d+1)^2$ and S -wise indep. rvars y_1, \dots, y_d over Ω each of which takes the vals 0,1 with prob. $1/2$.

The space and rvars. are explicitly constructed, given α a primitive elt. of field $F = GF(2^k)$ a k -dim'l algebra over $GF(2)$.

proof: [Everything instantiated for S , nothing special about $S!$].
~~MacWilliams and Sloane~~ (MacWilliams and Sloane)

Let $\alpha, \alpha^2, \dots, \alpha^d$ be the d non-zero elts of $GF(2^k)$ rep'd as col vectors of length k over $GF(2)$.

Let H be the parity check matrix of the extended BCH code of length d and distance $2t + 2 = 6$.

$$\begin{array}{c}
 \uparrow \\
 1 + 2 \cdot k \\
 \downarrow
 \end{array}
 \begin{array}{c}
 \leftarrow d \rightarrow \\
 \begin{pmatrix}
 1 & 1 & \dots & 1 \\
 \alpha & \alpha^2 & \dots & \alpha^d \\
 \alpha^3 & \alpha^6 & \dots & \alpha^{3d}
 \end{pmatrix}
 \end{array}
 \quad H \text{ is a matrix over } GF(2).$$

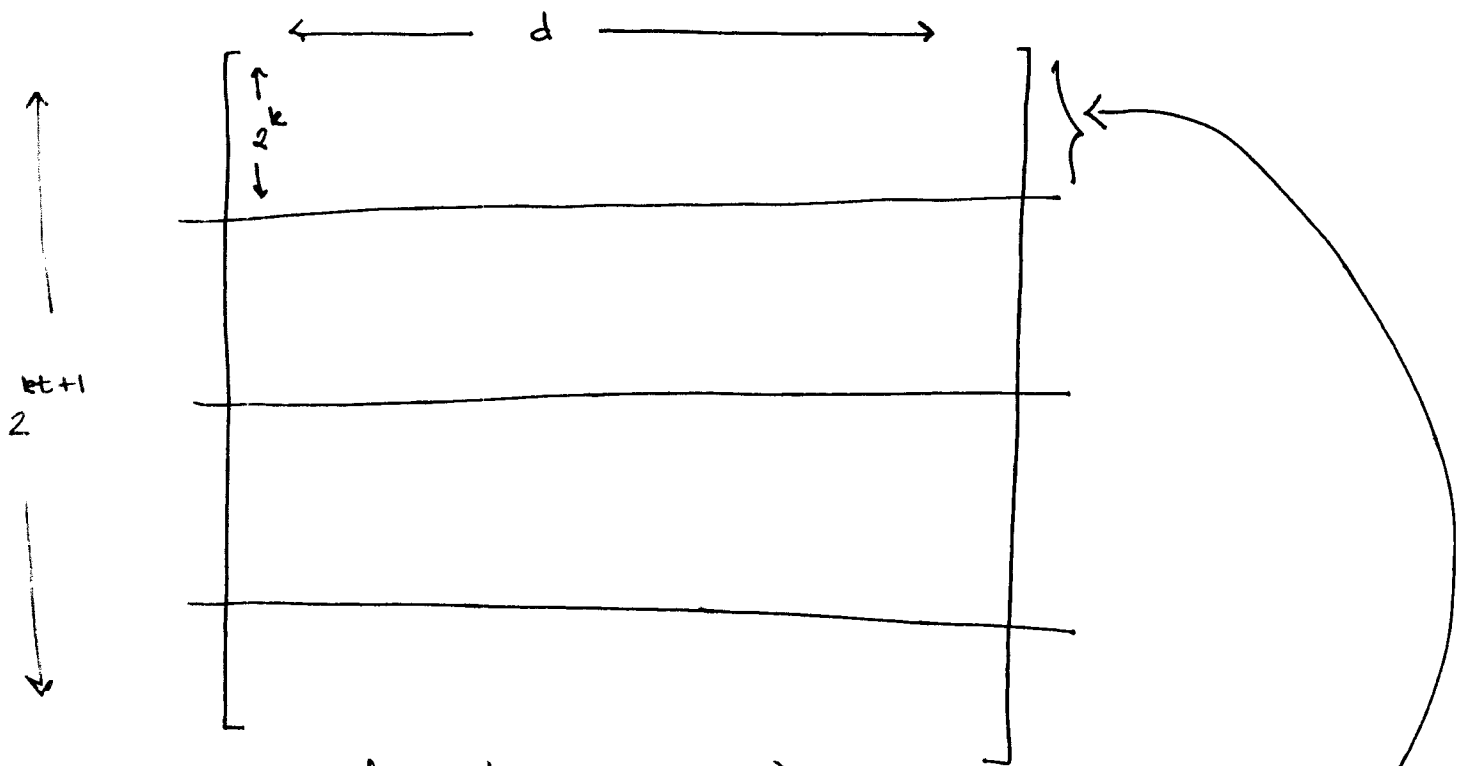
Lemma: Any set of $5 = 2 \cdot 2 + 1$ cols of H is linearly indep. over $GF(2)$.

proof: see MacWilliams & Sloane.

have a homog. system of $\overset{2t+1}{\uparrow}$ eq^s in $2t+1$ unknowns and matrix of coeffs. is a vandermonde matrix (which is nonsingular).

To generate rand vars., define $\Omega = \{1, 2, \dots, 2(d+1)^2\}$ and let $A = (a_{ij})$ $i \in \Omega$, $j = 1, \dots, d$ be a binary matrix with $2(d+1)^2 = 2^{kt+1}$ rows and d cols.

rows are all linear comb^s of the rows of H (over $GF(2)$).



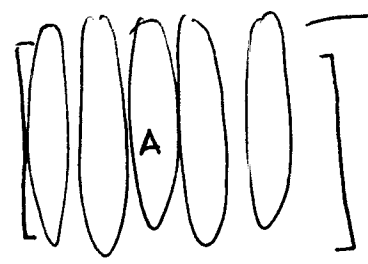
observe that for $b_i \in GF(2^k)$ $i=0,1,2$ we get one block of A

$$(b_0 \ b_1 \ b_2) \begin{pmatrix} 1 & 1 & \dots & 1 \\ \alpha & \alpha^2 & \dots & \alpha^d \\ \alpha^3 & \alpha^6 & \dots & \alpha^{3d} \end{pmatrix} = \left(\sum_{i=0}^2 b_i \alpha^i, \dots, \sum_{i=0}^2 b_i \alpha^{di} \right)$$

We endow the sample space Ω with the unif. prob. measure
Then y_j is defined by

$$y_j(i) = a_{ij} \text{ for } i \in \Omega \text{ and } j=1, \dots, d.$$

We need to show that y_j are 5-wise indep. and that each takes on 0,1 with equal prob.



Need to show that for every set J of up to 5 cols of A , the rows of submatrix

$$A_J = (a_{ij}) \text{ } i \in \Omega, j \in J$$

take on each of the 2^5 binary vectors equally often.

Lemma tells us that cols of submatrix H_J are linearly indep.

$$\Rightarrow \#(\text{rows of } A_J = \text{any given vector of length } 2^5)$$

$$= \#(\text{linear comb's of rows of } H_J = \text{vector})$$

$$= \#(\text{sols of system of 5 lin. indep. eq's in } kt+1 \text{ vars})$$

$$= 2^{kt+1-5} \quad //$$

To generate $y_j(i)$ quickly upon seeing j , pick i uniformly from $\Omega \iff$ pick setting of bit vectors (b_0, b_1, b_2) uniformly at random. There are $2 \cdot k + 1$ random bits here.
 $= 2 \cdot \log d + 1$

Then compute

$$\sum_{i=0}^{2^k-1} b_i \alpha^{j \cdot i} \in GF(2^k)$$

each of these are
vectors of length k

and α is stored
ahead of time.

⇒ Not storing A AND using a few random bits.