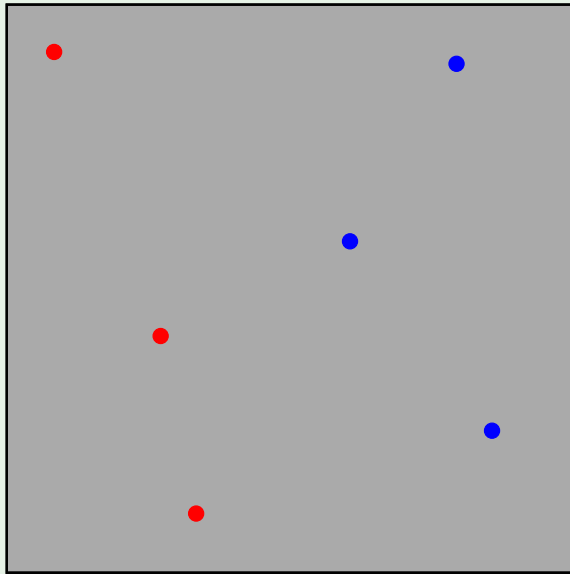


Review of Lecture 5

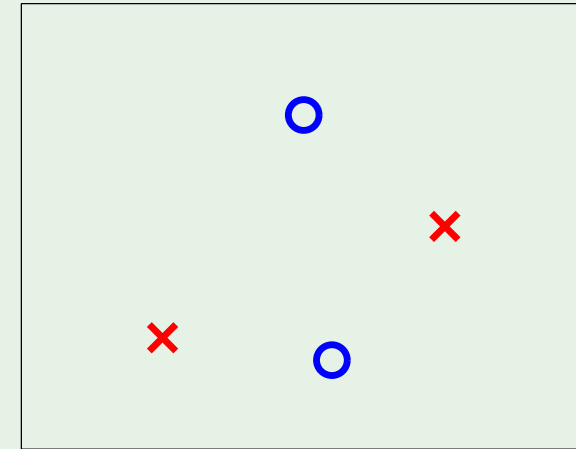
- Dichotomies = hypotheses restricted to a finite set of points



- Growth function

$$m_{\mathcal{H}}(N) = \max_{\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathcal{X}} |\mathcal{H}(\mathbf{x}_1, \dots, \mathbf{x}_N)|$$

- Break point



- Maximum # of dichotomies
resulting from the constraint of the break point (here k=2)

| \mathbf{x}_1 | \mathbf{x}_2 | \mathbf{x}_3 |
|----------------|----------------|----------------|
| ○ | ○ | ○ |
| ○ | ○ | ● |
| ○ | ● | ○ |
| ● | ○ | ○ |

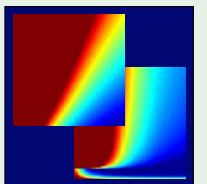
Learning From Data

Yaser S. Abu-Mostafa
California Institute of Technology

Lecture 6: Theory of Generalization



Sponsored by Caltech's Provost Office, E&AS Division, and IST • Thursday, April 19, 2012



Outline

- Proof that $m_{\mathcal{H}}(N)$ is polynomial (with a break point)
- Proof that $m_{\mathcal{H}}(N)$ can replace M

Bounding $m_{\mathcal{H}}(N)$

To show: $m_{\mathcal{H}}(N)$ is polynomial

We show: $m_{\mathcal{H}}(N) \leq \dots \leq \dots \leq$ a polynomial

We want to bound m , and we do this with $B(N, k)$ - the maximum number of dichotomies you can possibly have with no other constraints (other than N, k). This is purely combinatorial, meaning we can avoid any consideration of input space or correlation between events etc.

Key quantity:

$B(N, k)$: Maximum number of dichotomies on N points, with break point k

i.e max number of dichotomies on N points so that no k columns are shattered.

(non-specific to X or H , although H does determine k)

Recursive bound on $B(N, k)$

Consider the following table:

$$B(N, k) = \alpha + 2\beta$$

$B(N, k)$ is the maximum number of patterns we can get of N points such that no k columns have all possible patterns (are shattered).

S_1 contains rows which appear only once as far as x_1 to x_{N-1} are concerned - the prefix (x_1-x_{N-1}) happens once and only has one extension ($x_N=+1$ OR $x_N=-1$)

S_2 contains prefixes with both $x_N=+1$ AND $x_N=-1$ - we split each of these into subgroups S_2^+ and S_2^-

| | # of rows | \mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_{N-1} | \mathbf{x}_N |
|---------|-----------|----------------|----------------|----------|--------------------|----------------|
| S_1 | α | +1 | +1 | \dots | +1 | +1 |
| | | -1 | +1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | -1 | -1 |
| | | -1 | +1 | \dots | -1 | +1 |
| S_2^+ | β | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | +1 | +1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | -1 | +1 |
| S_2^- | β | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | -1 | -1 |

Estimating α and β

Focus on $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{N-1}$ columns:

$$\alpha + \beta \leq B(N-1, k)$$

All rows (total alpha+beta) highlighted are different, (note S_2^+ and S_2^- have equal prefixes, so not different).

Also, on the original matrix we could not find all possible patterns on any k columns, so we also cannot on the highlighted matrix. If we could, then these k columns would feature all possible patterns in the original matrix, but we do not. So the smaller matrix has the same break point k .

| | # of rows | \mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_{N-1} | \mathbf{x}_N |
|---------|-----------|----------------|----------------|----------|--------------------|----------------|
| S_1 | α | +1 | +1 | \dots | +1 | +1 |
| | | -1 | +1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | -1 | -1 |
| | | -1 | +1 | \dots | -1 | +1 |
| S_2^+ | β | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | +1 | +1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | -1 | +1 |
| S_2^- | β | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | -1 | -1 |

Estimating β by itself

Now, focus on the $S_2 = S_2^+ \cup S_2^-$ rows:

$$\beta \leq B(N-1, k-1)$$

| | # of rows | \mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_{N-1} | \mathbf{x}_N |
|--------------------|--------------------|----------------|----------------|----------|--------------------|----------------|
| S_1 | α | +1 | +1 | \dots | +1 | +1 |
| | | -1 | +1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | -1 | -1 |
| | | -1 | +1 | \dots | -1 | +1 |
| S_2 | S_2^+ β | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | +1 | +1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | +1 |
| | | -1 | -1 | \dots | -1 | +1 |
| S_2^- β | β | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | +1 | -1 |
| | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | +1 | -1 | \dots | +1 | -1 |
| | | -1 | -1 | \dots | -1 | -1 |

Putting it together

$$B(N, k) = \alpha + 2\beta$$

$$\alpha + \beta \leq B(N - 1, k)$$

$$\beta \leq B(N - 1, k - 1)$$

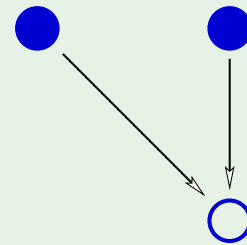
$$B(N, k) \leq$$

$$B(N - 1, k) + B(N - 1, k - 1)$$

| | | # of rows | \mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_{N-1} | \mathbf{x}_N |
|-------|----------|-----------|----------------|----------------|----------|--------------------|----------------|
| S_1 | α | | +1 | +1 | \dots | +1 | +1 |
| | | | -1 | +1 | \dots | +1 | -1 |
| | | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | | +1 | -1 | \dots | -1 | -1 |
| | | | -1 | +1 | \dots | -1 | +1 |
| S_2 | S_2^+ | β | +1 | -1 | \dots | +1 | +1 |
| | | | -1 | -1 | \dots | +1 | +1 |
| | | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | | +1 | -1 | \dots | +1 | +1 |
| | | | -1 | -1 | \dots | -1 | +1 |
| | S_2^- | β | +1 | -1 | \dots | +1 | -1 |
| | | | -1 | -1 | \dots | +1 | -1 |
| | | | \vdots | \vdots | \vdots | \vdots | \vdots |
| | | | +1 | -1 | \dots | +1 | -1 |
| | | | -1 | -1 | \dots | -1 | -1 |

Numerical computation of $B(N, k)$ bound

$$B(N, k) \leq B(N - 1, k) + B(N - 1, k - 1)$$



| | | k | | | | | | |
|-----|---|-----|---|----|----|----|----|----|
| | | 1 | 2 | 3 | 4 | 5 | 6 | .. |
| N | 1 | 1 | 2 | 2 | 2 | 2 | 2 | .. |
| | 2 | 1 | 3 | 4 | 4 | 4 | 4 | .. |
| | 3 | 1 | 4 | 7 | 8 | 8 | 8 | .. |
| | 4 | 1 | 5 | 11 | .. | .. | .. | .. |
| | 5 | 1 | 6 | : | . | | | |
| | 6 | 1 | 7 | : | | . | | |
| | : | : | : | : | | | . | |

Analytic solution for $B(N, k)$ bound

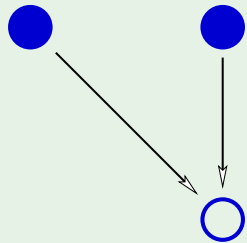
$$B(N, k) \leq B(N - 1, k) + B(N - 1, k - 1)$$

Theorem:

$$B(N, k) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

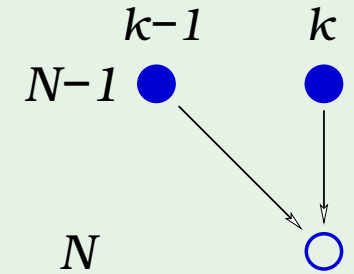
1. Boundary conditions: easy

| | | k | | | | | | |
|-----|---|-----|---|---|---|---|---|----|
| | | 1 | 2 | 3 | 4 | 5 | 6 | .. |
| N | 1 | 1 | 2 | 2 | 2 | 2 | 2 | .. |
| | 2 | 1 | | | | | | |
| | 3 | 1 | | | | | | |
| | 4 | 1 | | | | | | |
| | 5 | 1 | | | | | | |
| | 6 | 1 | | | | | | |
| | : | : | | | | | | |



2. The induction step

$$\begin{aligned}
 \sum_{i=0}^{k-1} \binom{N}{i} &= \sum_{i=0}^{k-1} \binom{N-1}{i} + \sum_{i=0}^{k-2} \binom{N-1}{i} \text{ ?} \\
 &= 1 + \sum_{i=1}^{k-1} \binom{N-1}{i} + \sum_{i=1}^{k-1} \binom{N-1}{i-1} \\
 &= 1 + \sum_{i=1}^{k-1} \left[\binom{N-1}{i} + \binom{N-1}{i-1} \right] \\
 &= 1 + \sum_{i=1}^{k-1} \binom{N}{i} = \sum_{i=0}^{k-1} \binom{N}{i} \checkmark
 \end{aligned}$$



It is polynomial!

For a given \mathcal{H} , the break point k is fixed

$$m_{\mathcal{H}}(N) \leq \underbrace{\sum_{i=0}^{k-1} \binom{N}{i}}_{\text{maximum power is } N^{k-1}}$$

Three examples

$$\sum_{i=0}^{k-1} \binom{N}{i}$$

- \mathcal{H} is positive rays: (break point $k = 2$)

$$m_{\mathcal{H}}(N) = N + 1 \leq N + 1$$

- \mathcal{H} is positive intervals: (break point $k = 3$)

$$m_{\mathcal{H}}(N) = \frac{1}{2}N^2 + \frac{1}{2}N + 1 \leq \frac{1}{2}N^2 + \frac{1}{2}N + 1$$

- \mathcal{H} is 2D perceptrons: (break point $k = 4$)

$$m_{\mathcal{H}}(N) = ? \leq \frac{1}{6}N^3 + \frac{5}{6}N + 1$$

Outline

- Proof that $m_{\mathcal{H}}(N)$ is polynomial
- Proof that $m_{\mathcal{H}}(N)$ can replace M

What we want

Instead of:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2 \quad \textcolor{red}{M} \quad e^{-2\epsilon^2 N}$$

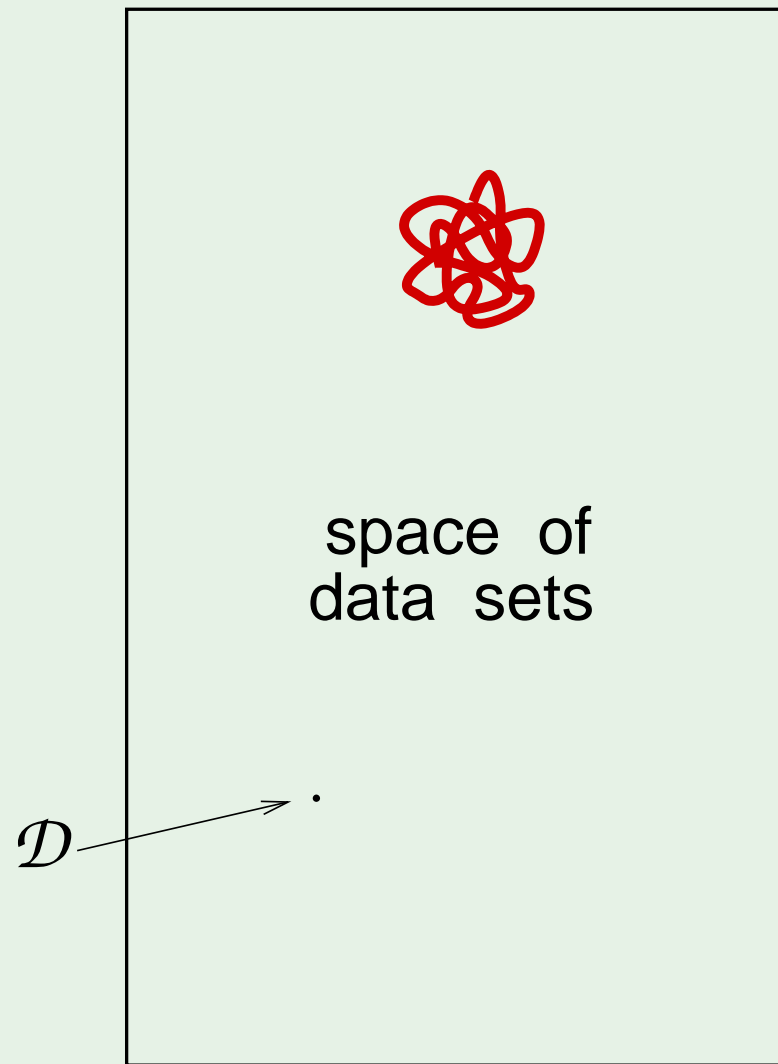
We want:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2 \quad \textcolor{red}{m}_{\mathcal{H}}(N) \quad e^{-2\epsilon^2 N}$$

Pictorial proof ☺

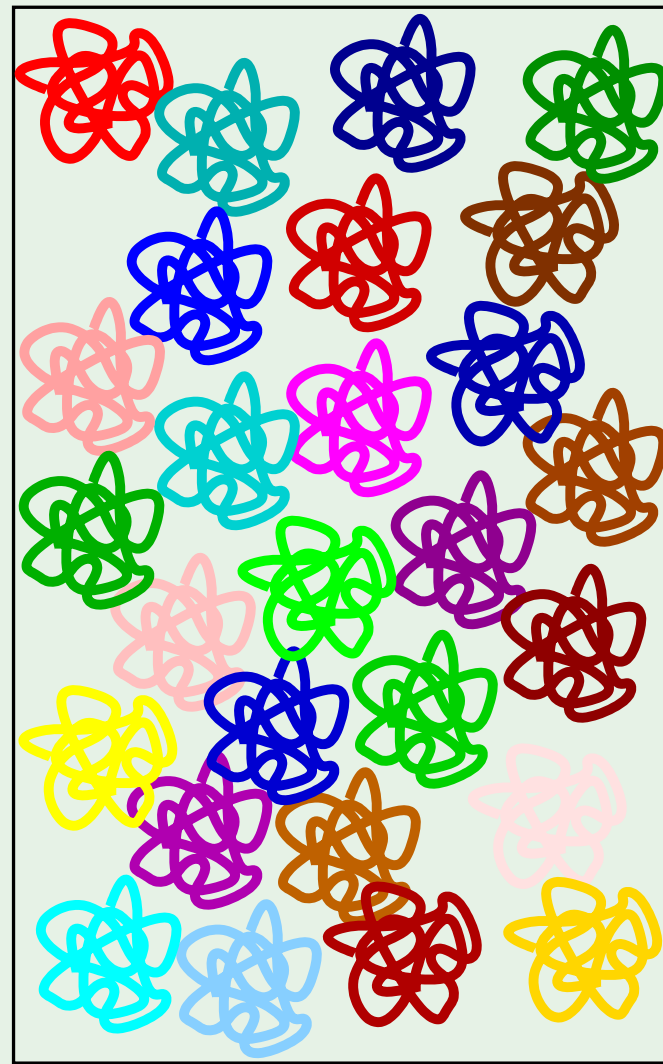
- How does $m_{\mathcal{H}}(N)$ relate to overlaps? since M is created from the union bound which assumes disjoint hypotheses
- What to do about E_{out} ? since the growth function relies on a finite sample (and the subsequent dichotomies), so it will handle the E_{in} aspect of Hoeffding. However E_{out} relates to the performance over the entire input space X and so we are dealing with full hypotheses, not dichotomies, so we lose the benefit of the growth function m .
- Putting it together

Hoeffding Inequality



(a)

Union Bound



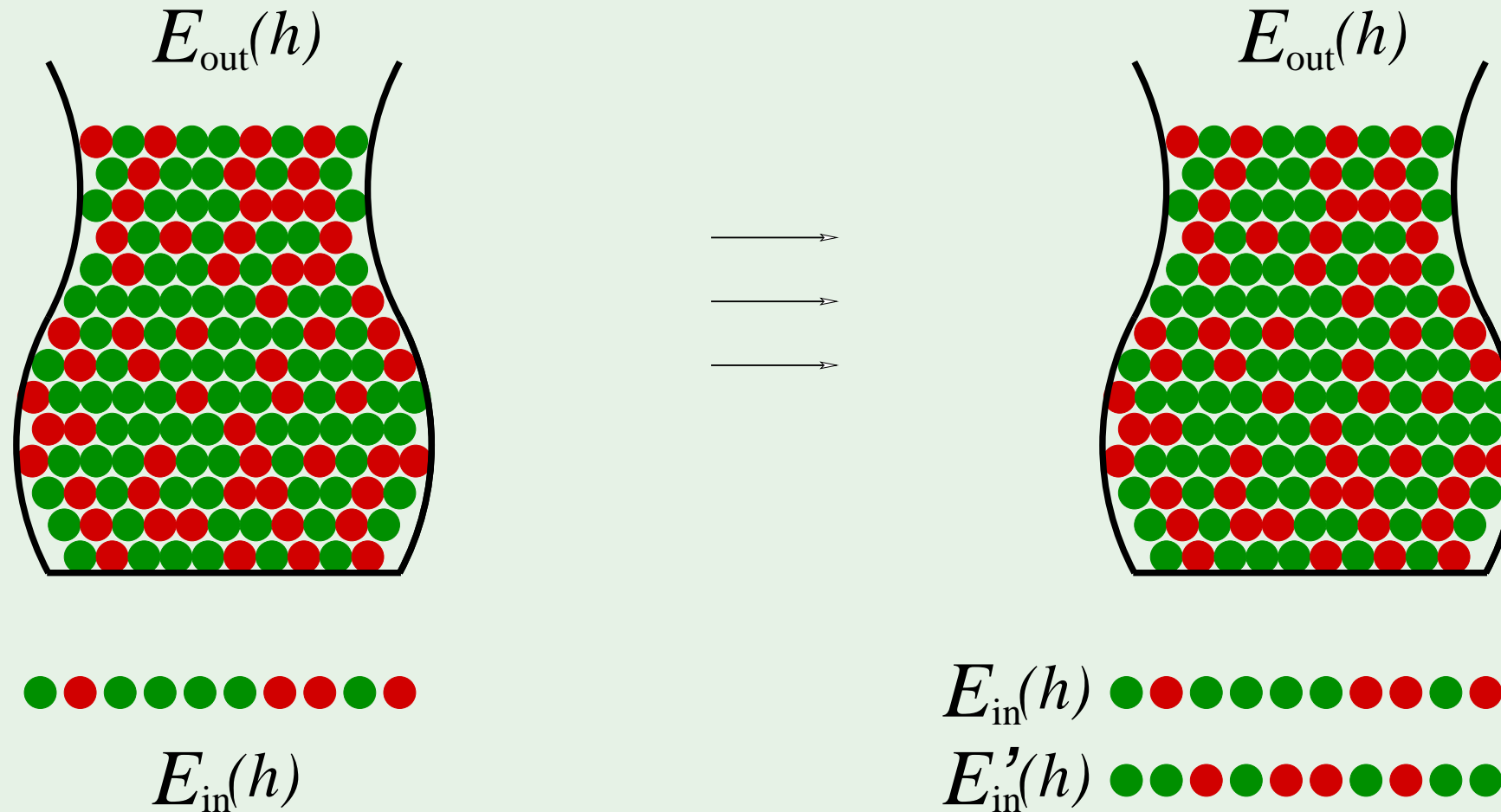
(b)

VC Bound



(c)

What to do about E_{out}



E_{in} and E'_{in} track each other since they both track E_{out} (even if their tracking is looser) - e.g. you expect two polls of equal N to have close results to each other. Like how the tracking of E_{out} and E_{in} become looser as the number of hypotheses increased (from M in Hoeffding), it also happens with E_{in} and E'_{in} . If we characterize this using the two samples only, no longer appealing to E_{out} , we are completely in the realm of dichotomies and, although the sample is bigger ($2N$), we can define a growth function on them - see next slide.

Putting it together

Not quite:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2 m_{\mathcal{H}}(N) e^{-2 \epsilon^2 N}$$

but rather:

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 4 m_{\mathcal{H}}(2N) e^{-\frac{1}{8} \epsilon^2 N}$$

The Vapnik-Chervonenkis Inequality