

Lecture 3

September 13, 2023

Instructor: Sepehr Assadi

Disclaimer: These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.

Topics of this Lecture

1	Threshold Testing Distinct Elements	1
1.1	The Algorithm	1
1.2	Intuition	2
1.3	Formal Analysis	2
2	Independence for Space: Limited-Independence Hash Functions	4
2.1	Improving Space Complexity of Algorithm 1	5
3	The Original Distinct Elements Problem?	6

1 Threshold Testing Distinct Elements

We finished the last lecture with the following problem.

Problem 1. At the beginning of the stream, you are given a value $\widetilde{\text{DE}} \in [m]$ and a parameter $\varepsilon \in (0, 1)$. Then, you are given a stream of n numbers (possibly with repetitions) from a universe $[m]$; we denote the number of distinct elements in the stream by DE . The goal is to, with probability at least $2/3$, output *Yes* if $\text{DE} \geq \widetilde{\text{DE}}$ and output *No* if $\text{DE} < (1 - \varepsilon) \cdot \widetilde{\text{DE}}$; if the value of DE is between these two numbers, either answer is considered correct.

And, we promised to give an algorithm that solves this problem using $\text{poly}(\log n, \log m, 1/\varepsilon)$ bits of space. We will present such an algorithm in this lecture.

Simplifying assumption. Throughout this lecture, without loss of generality, we assume that $\widetilde{\text{DE}} \geq 100/\varepsilon^2$. This is because otherwise, we can simply use the deterministic $O(\widetilde{\text{DE}} \cdot \log m)$ -space naive algorithm that stores all the distinct elements it sees and answers *Yes* as soon as it sees $T+1$ of them; when $\widetilde{\text{DE}} < 100/\varepsilon^2$, this algorithm only requires $O(\log m/\varepsilon^2)$ space which is sufficient for our purpose.

1.1 The Algorithm

The following algorithm solves **Problem 1**, for a given threshold $\widetilde{\text{DE}}$.

$$\frac{100}{\varepsilon^2} \log m$$

$$\frac{\log m}{\varepsilon^2}$$

Algorithm 1. An algorithm for threshold testing distinct elements for a given $\widetilde{DE} \in [m]$ and $\varepsilon \in (0, 1)$:

- (i) Let $t := 12/\varepsilon^2$ and pick a **hash function** $h : [m] \rightarrow [\widetilde{DE}/t]$ uniformly at random from the set of all functions from $[m] \rightarrow [\widetilde{DE}/t]$.
- (ii) Count the number of *distinct* elements in the stream that are hashed to 1, i.e., count the size of the set $\{e \mid h(e) = 1\}$, denoted by X .
- (iii) Return *Yes* if the number X calculated above is at least $(1 - \varepsilon/2) \cdot t$ or *No* otherwise.

We will talk about the space complexity of this algorithm later in the lecture. For now, we should only note that *as it is*, this algorithm requires prohibitively large space to store the hash function h and hence is not space-efficient. Although the rest of the algorithm uses $O(t \cdot \log m) = O(\log m / \varepsilon^2)$ bits of space which is good enough for us. Regardless, almost all the main ideas of the algorithm are already here and thus we focus on proving its correctness in the following.

1.2 Intuition

$$\frac{1}{\frac{\widetilde{DE}}{t}}$$

$$O\left(\left(1 - \frac{\varepsilon}{2}\right) \cdot \frac{12}{\varepsilon^2} \cdot \log m\right) = O\left(\left(\frac{12}{\varepsilon^2} - \frac{6}{\varepsilon}\right) \cdot \log m\right) = O\left(\frac{\log m}{\varepsilon^2}\right)$$

The intuition behind the algorithm is as follows: for each *distinct* number in the stream, $h(\cdot)$ has a chance of hitting 1 with probability t/\widetilde{DE} . As such, if $DE \geq \widetilde{DE}$, then it is very likely that a “good number” of the element will be hashed to 1, but if $DE \leq (1 - \varepsilon) \cdot \widetilde{DE}$ that number should be considerably lower (both cases in a probabilistic sense).

You may ask what is the role of t then, i.e., why did we pick t to be $\Omega(1/\varepsilon^2)$ and not simply, say, 1?¹ This is done for the purpose of “*variance reduction*”: see the calculation for the variance of the random variables and how it is used in the analysis to see the necessity of using a larger t .

1.3 Formal Analysis

We now formalize this intuition and prove the correctness of the algorithm.

Lemma 1. On any input stream and for any choice of parameters $\varepsilon \in (0, 1)$ and $\widetilde{DE} \geq 100/\varepsilon^2$:

- If $DE \geq \widetilde{DE}$, then **Algorithm 1** outputs *Yes* with probability at least $2/3$;
- If $DE < (1 - \varepsilon) \cdot \widetilde{DE}$, then **Algorithm 1** outputs *No* with probability at least $2/3$.

Proof. Let $\{e_1, \dots, e_{DE}\}$ denote the distinct elements in the stream. For $i \in [DE]$, define the indicator random variable $X_i \in \{0, 1\}$ which is 1 iff $h(e_i) = 1$. Under this definition, the random variable X in **Algorithm 1** is:

$$X = \sum_{i=1}^{DE} X_i.$$

We can thus calculate the expected value of X as follows:

$$\begin{aligned} \mathbb{E}[X] &= \mathbb{E}\left[\sum_{i=1}^{DE} X_i\right] && \text{(by the equation above)} \\ &= \sum_{i=1}^{DE} \mathbb{E}[X_i] && \text{(by linearity of expectation)} \end{aligned}$$

¹Notice that we ideally want t to be as small as possible as the space of the algorithm depends linearly on t .

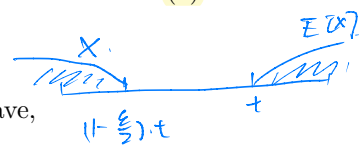
$$\begin{aligned}
&= \sum_{i=1}^{\text{DE}} \Pr(h(e_i) = 1) && \text{(by the definition of indicator } X_i) \\
&= \sum_{i=1}^{\text{DE}} \frac{t}{\widetilde{\text{DE}}} && \text{(as } h(\cdot) \text{ maps each element to a uniformly random position in } [\widetilde{\text{DE}}/t]) \\
&= \text{DE} \cdot \frac{t}{\widetilde{\text{DE}}}. && (1)
\end{aligned}$$

For our analysis, we also need to bound the variance of X , which can be done as follows:

$$\begin{aligned}
\text{Var}[X] &= \text{Var}\left[\sum_{i=1}^{\text{DE}} X_i\right] && \text{(again, by the equation } X = \sum_i X_i) \\
&= \sum_{i=1}^{\text{DE}} \text{Var}[X_i] && \text{(variance of sum of independent variables is sum of variances)} \\
&\leq \sum_{i=1}^{\text{DE}} \mathbb{E}[X_i^2] && \text{(by the definition of variance } \text{Var}[X_i] = \mathbb{E}[X_i^2] - (\mathbb{E}[X_i])^2) \\
&= \sum_{i=1}^{\text{DE}} \mathbb{E}[X_i] && \text{(as } X_i^2 = X_i \text{ since } X_i \in \{0, 1\}) \\
&= \mathbb{E}[X]. && (2)
\end{aligned}$$

We can now analyze each case separately.

Case I: when $\text{DE} \geq \widetilde{\text{DE}}$: In this case, by the bound on the expectation in Eq (1), we have,



$$X < (1 - \frac{\epsilon}{2}) \cdot \mathbb{E}[X] \implies \mathbb{E}[X] \geq t.$$

$$|\mathbb{E}[X] - X| \geq |\mathbb{E}[X] - (1 - \frac{\epsilon}{2}) \cdot \mathbb{E}[X]| = \frac{\epsilon}{2} \cdot \mathbb{E}[X]$$

why in $\mathbb{E}[X]$?
next page.

The algorithm will say *No* if $X < (1 - \epsilon/2) \cdot t$; in this case, this requires X to deviate by at least $(\epsilon/2) \cdot \mathbb{E}[X]$ from its expectation. This is exactly the topic of concentration inequalities. In particular, recall Chebyshev's inequality from the last lecture:

Proposition 2 (Chebyshev's Inequality). For any random variable X and $b > 0$,

$$\Pr(|X - \mathbb{E}[X]| \geq b) \leq \frac{\text{Var}[X]}{b^2}.$$

We can apply Chebyshev's inequality to our random variable X to get:

$$\Pr(\text{Algorithm 1 says No in Case I}) \leq \Pr(|X - \mathbb{E}[X]| \geq (\epsilon/2) \cdot \mathbb{E}[X]) \quad \text{(as described above)}$$

$$\frac{\text{Var}[X]}{\frac{\epsilon^2}{4} \cdot \mathbb{E}[X]^2}$$

$$\leq \frac{4 \cdot \text{Var}[X]}{\epsilon^2 \cdot \mathbb{E}[X]^2}$$

(by Chebyshev's inequality of Proposition 2 for $b = (\epsilon/2) \cdot \mathbb{E}[X]$)

$$\leq \frac{4}{\epsilon^2 \cdot \mathbb{E}[X]} \quad \text{(by the upper bound of } \mathbb{E}[X] \text{ on variance in Eq (2))}$$

$$\leq \frac{4}{\epsilon^2 \cdot (12/\epsilon^2)}$$

(by the lower bound of t on expectation and the choice of $t = 12/\epsilon^2$)

$$= \frac{1}{3}.$$

This proves the first bullet of the lemma statement.

$$\mathbb{E}[X] = D\tilde{E} \cdot \frac{t}{D\tilde{E}}$$

Case II: when $DE < (1 - \varepsilon) \cdot \tilde{DE}$: In this case, by the bound on the expectation in Eq (1), we have,

$$(1 - \varepsilon/2) \cdot t \quad (1 - \varepsilon/2) \cdot t$$

$$\mathbb{E}[X] < (1 - \varepsilon) \cdot t.$$

$$|X - \mathbb{E}[X]| \geq \left| (1 - \frac{\varepsilon}{2}) \cdot t - (1 - \varepsilon) \cdot t \right| = \frac{\varepsilon}{2} t.$$

The algorithm will say *Yes* if $X \geq (1 - \varepsilon/2) \cdot t$; in this case, this requires X to deviate by at least $(\varepsilon/2) \cdot t$ from its expectation.² We thus have,

$$\begin{aligned} \Pr(\text{Algorithm 1 says Yes in Case II}) &\leq \Pr(|X - \mathbb{E}[X]| > (\varepsilon/2) \cdot t) && \text{(as described above)} \\ &\leq \frac{4 \cdot \text{Var}[X]}{\varepsilon^2 \cdot t^2} \\ &\quad \text{(by Chebyshev's inequality of Proposition 2 for } b = (\varepsilon/2) \cdot \mathbb{E}[X]) \\ &\leq \frac{4 \cdot \mathbb{E}[X]}{\varepsilon^2 \cdot t^2} && \text{(by the upper bound of } \mathbb{E}[X] \text{ on variance in Eq (2))} \\ &\leq \frac{4 \cdot (1 - \varepsilon)}{\varepsilon^2 \cdot (12/\varepsilon^2)} \\ &\quad \text{(by the upper bound of } \mathbb{E}[X] < (1 - \varepsilon) \cdot t \text{ and the choice of } t = 12/\varepsilon^2) \\ &< \frac{1}{3}. \end{aligned}$$

This proves the second bullet of the lemma statement and concludes the whole proof. \square

Thus, the algorithm has a probability of success of at least $2/3$ for our problem, as desired. In the next part, we see how to fix the issue with the space complexity of the problem—in particular, how to replace the hash function h with something that we can store more efficiently.

2 Independence for Space: Limited-Independence Hash Functions

Generating and storing a random function $h : [a] \rightarrow [b]$ requires $O(a \log b)$ bits, which is often too costly. In the case of Algorithm 1 this amounts to $\omega(m)$ space which makes the whole algorithm entirely useless! Fortunately however, we can make the analysis of Algorithm 1 work even if the hash function h we picked is not *completely random*, but has some *limited independence* only. Let us define this formally as follows.

Definition 3. A family $\mathcal{H} = \{h : [a] \rightarrow [b]\}$ is called a **k -wise independent** family of hash functions if for all pairwise distinct $x_1, \dots, x_k \in [a]$ and all $y_1, \dots, y_k \in [b]$,

$$\Pr_{h \sim \mathcal{H}}(h(x_1) = y_1 \wedge \dots \wedge h(x_k) = y_k) = \frac{1}{b^k}.$$

\uparrow
' \sim ' means h is u.a.r. chosen from \mathcal{H}

Observe that a k -wise independent family may also be $(k+1)$ -wise independent, i.e., the definition does not necessarily break for $k+1$ hash values (although for “interesting” families this is almost always the case). For instance, a truly random hash function is k -wise independent for all $k \in [m]$.

Proposition 4. Let $\mathcal{H} = \{h : [a] \rightarrow [b]\}$ be a k -wise independent, and $h \sim \mathcal{H}$ chosen at random. Let $x_1, \dots, x_k \in [a]$ be arbitrary pairwise distinct elements. Then:

1. for every $i \in [k]$, $h(x_i)$ is uniform over $[b]$;
2. $h(x_1), \dots, h(x_k)$ are mutually independent.

²Notice that this time, we bound the deviation as a function of t and *not* $\mathbb{E}[X]$; this is because, in this case, it is possible for $\mathbb{E}[X]$ to be very small and thus bounding the deviation just as a function of expectation will be too weak. You are also encouraged to check that in the previous case, bounding the deviation as a function of t does *not* work (because $\mathbb{E}[X]$ can be much larger than t in that case).

Proof. We prove each part separately.

1. We only prove this for $i = 1$; the rest is symmetric. Let $y_1 \in [b]$. Observe that

$$\begin{aligned} \Pr_{h \sim \mathcal{H}}(h(x_1) = y_1) &= \sum_{y_2, \dots, y_k \in [b]} \Pr_{h \sim \mathcal{H}}(h(x_1) = y_1 \wedge h(x_2) = y_2 \wedge \dots \wedge h(x_k) = y_k) \\ &\quad \text{(partitioning the sample space)} \\ &= \sum_{y_2, \dots, y_k \in [b]} \frac{1}{b^k} = \frac{b^{k-1}}{b^k} = \frac{1}{b}. \end{aligned} \quad \text{(by definition of } k\text{-wise independent family)}$$

2. Let $y_1, \dots, y_k \in [b]$. Since all $h(x_i)$ are uniform over $[b]$, it follows that

$$\Pr_{h \sim \mathcal{H}}(h(x_1) = y_1 \wedge \dots \wedge h(x_k) = y_k) = \frac{1}{b^k} = \prod_{i=1}^k \Pr_{h \sim \mathcal{H}}(h(x_i) = y_i).$$

This concludes the proof. □

Example. Let $k \geq 2$ be an integer and $p > k$ be a prime number. Here is an example of a k -wise independent family of hash functions mapping $[p] \rightarrow [p]$

A k -wise Independent Family of Hash Functions on $[p] \rightarrow [p]$ for a prime p :

1. \mathcal{H} is the set of degree- $(k-1)$ polynomials over \mathbb{F}_p (field of integers mod prime p). That is,

$$\mathcal{H} := \{h : [p] \rightarrow [p] \mid h(x) = c_{k-1} \cdot x^{k-1} + c_{k-2} \cdot x^{k-2} + \dots + c_1 \cdot x + c_0, \text{ with } c_0, \dots, c_{k-1} \in \mathbb{F}_p\}.$$

2. Sample $c_0, \dots, c_{k-1} \in \mathbb{F}_p$ and return h as the polynomial defined by these coefficients.

To see why this is a k -wise independent hash function, note that any degree- $(k-1)$ polynomial h is uniquely determined by having k of its values (i.e., k distinct $(x, h(x))$ pairs for $x \in \mathbb{F}_p$): if we fix only $k-1$ values of h on x_1, \dots, x_{k-1} , value of $h(x_k)$ for any other x_k is still chosen uniformly at random from \mathbb{F}_p (we omit the simple algebraic proof of this statement as it is not the focus of this lecture/course).

The important thing we would like to note about the family \mathcal{H} is on how much space we need to store h . Since each function in the family is defined by k polynomial coefficients, the space required to generate and store it is only $O(k \log p)$ bits (as opposed to $O(p \log p)$ for a truly random hash function mapping $[p] \rightarrow [p]$). It is also possible to evaluate any such hash function in the same amount of space.

Although this family only works for $a = b = p$, we can in general construct families for arbitrary a and b .

Proposition 5. For any integers $a, b \geq 1$, there exists a k -wise independent family of hash functions mapping $[a] \rightarrow [b]$, that requires $O(k \cdot (\log a + \log b))$ bits.

2.1 Improving Space Complexity of Algorithm 1

To make Algorithm 1 space-efficient, we are going to replace the truly random hash function $h : [m] \rightarrow [\widetilde{DE}/t]$ with a $pair$ -wise independent hash function instead. By Proposition 5, this means we can store h in only $O(\log m)$ bits (as $\widetilde{DE}/t \leq m$) and evaluate $h(\cdot)$ on each arriving element quickly and in a space-efficient manner still. The rest of the algorithm also needed $O(\log m/\varepsilon^2)$ bits, thus we now have a truly space-efficient algorithm for the problem.

But, what about the analysis? There were only two key places that we used the properties h :

- In Eq (1) to compute the expected value of X . In particular, we used the fact that for each element e in the universe, $\Pr(h(e_i) = 1) = t/\widetilde{DE}$, namely, that $h(e_i)$ is distributed uniformly over the range of h . This continues to hold for any pair-wise independent hash function also as argued in Proposition 4.

As an aside, this is an easy property to satisfy that holds even for “weaker” hash families, say *universal* hash families, or even *uniform* hash families—for instance, consider the family of functions $h : [a] \rightarrow [b]$ consisting of b functions $\{h_i \mid h_i(x) = i \ \forall x \in [a]\}$; this is obviously a “bad” hash family that maps *all* the elements to the same position, and still even this family is enough to satisfy the property.

- In Eq (2) to say that variance of the sum is equal to sum of the variances because X_i ’s are chosen independently. We no longer have the independence property here. However, the conclusion, namely, that variance of the sum is equal to the sum of variances holds even for pairwise independent variables—which X_i ’s are because they are deterministic functions $h(e_i)$ ’s, which are pairwise independent—as we prove below.

Proposition 6. *Let X_1, \dots, X_n be a family of n pair-wise independent variables. Then,*

$$\text{Var} \left[\sum_{i=1}^n X_i \right] = \sum_{i=1}^n \text{Var} [X_i].$$

Proof. We have

$$\begin{aligned} \text{Var} \left[\sum_{i=1}^n X_i \right] &= \mathbb{E} \left[\left(\sum_{i=1}^n X_i \right)^2 \right] - \left(\mathbb{E} \left[\sum_{i=1}^n X_i \right] \right)^2 && \text{(by the definition of variance)} \\ &= \mathbb{E} \left[\sum_{i=1}^n X_i^2 + \sum_{i \neq j} X_i \cdot X_j \right] - \left[\sum_{i=1}^n \mathbb{E} [X_i]^2 - \sum_{i \neq j} \mathbb{E} [X_i] \cdot \mathbb{E} [X_j] \right] && \text{(by expanding the sums)} \\ &= \left(\sum_{i=1}^n \left(\mathbb{E} [X_i^2] - \mathbb{E} [X_i]^2 \right) \right) + \left(\sum_{i \neq j} \left(\mathbb{E} [X_i \cdot X_j] - \mathbb{E} [X_i] \cdot \mathbb{E} [X_j] \right) \right) && \text{(by linearity of expectation and re-ordering the terms)} \\ &= \sum_{i=1}^n \text{Var} [X_i] + 0, && \text{=E[Xi] E[Xj]} \end{aligned}$$

(first term by definition of variance, second term because $X_i \perp X_j$ for pair-wise independent variables)

which proves the result. \square

This means that the modification of Algorithm 1, by replacing h with a pairwise independent hash functions, works exactly as before and thus solves our problem, but now with a much better space. This proves the following theorem.

Theorem 7. *There is a streaming algorithm for threshold testing number of distinct elements in Problem 1 that uses $O(\log m/\varepsilon^2)$ space and outputs the correct answer with probability at least $2/3$.*

3 The Original Distinct Elements Problem?

What about the original problem: estimating the number of distinct elements instead of threshold testing?

There is a simple way to go from Theorem 7 to that problem: run the algorithm of Theorem 7 in parallel for different thresholds $\widetilde{DE} \in \{1, (1+\varepsilon), (1+\varepsilon)^2, (1+\varepsilon)^3, \dots, m\}$; then, find the *largest* choice of \widetilde{DE} for which the algorithm returns *Yes*, and output that as the estimate of DE for the stream. Notice that this

increasing the space by a factor of $O(\log_{(1+\epsilon)}(m)) = O((\log m)/\epsilon)$ which is okay for our purpose. Assuming every application of [Theorem 7](#) in this process is also correct, it is easy to see that the returned answer will be a $(1 + \Theta(\epsilon))$ -approximation of DE – we can then use a smaller ϵ in the algorithm, if needed, by changing the space with a constant factor, and obtain a truly $(1 + \epsilon)$ -approximation.

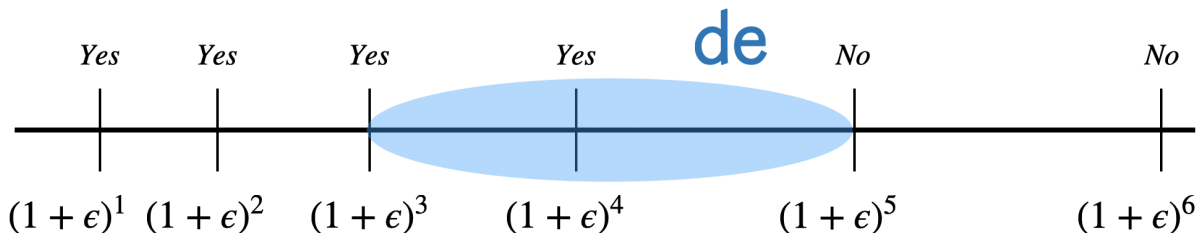


Figure 1: An illustration of the reduction from original estimation problem to threshold testing. Assuming all *Yes/No* responses are correct, the true estimate should lie somewhere in the marked (blue) region—the reason for the gap is to account for the fact that in the threshold testing, the answer can be arbitrary when neither case are satisfied.

An important caveat is that, as stated, [Theorem 7](#) only guarantees $2/3$ probability of success and thus it is not going to be the case that all of its $O((\log m)/\epsilon)$ invocations in this strategy return a correct answer. As such, we first need to *boost* the probability of success of the algorithm to a larger value before running this strategy—this will be the content of a future lecture.

Finally, we note that the above strategy of going from threshold testing to the original estimation problem, as well as boosting success of randomized algorithms, is a completely generic idea that has nothing to do with the distinct element problem we studied, and can be applied to most other settings and problems.

Remark. The distinct element problem—in this formulation—was first studied by Flajolet and Martin in [\[FM83\]](#), long before the formalization of the streaming model, as was revisited in the work of Alon, Matias, and Szegedy [\[AMS96\]](#) that pioneered the streaming model.

The algorithm we discussed in this lecture was inspired by an algorithm of Bar-Yossef, Jayram, Kumar, Sivakumar, and Trevisan [\[BJK⁺02\]](#). This algorithm was later improved in a series of work, culminating in the asymptotically optimal algorithm of Kane, Nelson, and Woodruff [\[KNW10\]](#) with space complexity $O(\frac{1}{\epsilon^2} + \log m)$ bits.

References

- [AMS96] Noga Alon, Yossi Matias, and Mario Szegedy. The space complexity of approximating the frequency moments. In *Proceedings of the Twenty-Eighth Annual ACM Symposium on the Theory of Computing, Philadelphia, Pennsylvania, USA, May 22-24, 1996*, pages 20–29, 1996. [7](#)
- [BJK⁺02] Ziv Bar-Yossef, T. S. Jayram, Ravi Kumar, D. Sivakumar, and Luca Trevisan. Counting distinct elements in a data stream. In *Randomization and Approximation Techniques, 6th International Workshop, RANDOM 2002, MA, USA, September 13-15, 2002*, pages 1–10, 2002. [7](#)
- [FM83] Philippe Flajolet and G. Nigel Martin. Probabilistic counting. In *24th Annual Symposium on Foundations of Computer Science, Arizona, USA, 7-9 November 1983*, pages 76–82, 1983. [7](#)
- [KNW10] Daniel M. Kane, Jelani Nelson, and David P. Woodruff. An optimal algorithm for the distinct elements problem. In *Proceedings of the Twenty-Ninth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems, PODS 2010, June 6-11, 2010, Indianapolis, Indiana, USA*, pages 41–52, 2010. [7](#)