

## 6.7720/18.619/15.070 Lecture 3

# The First & Second Moment Methods: Statistical Inference for the Broadcast Process

Kuikui Liu

February 10, 2025

**Acknowledgements & Disclaimers** *In the process of writing these notes, we consulted materials created by Guy Bresler and David Gamarnik, who taught previous iterations of this course. We are grateful for the discussions we had with them. We also consulted materials by Sebastien Roch, as well as the paper [Bor+06]. Please be advised that these notes have not been subjected to the usual scrutiny reserved for formal publications. If you do spot an error, please contact the instructor.*

## 1 Broadcasting on Trees

A fundamental problem in statistical inference is that of *hypothesis testing*. We are given a *sample* (often called the “observation”)  $X$  from some mystery probability distribution, and there are two possibilities (or *hypotheses*)  $\mu, \nu$  for what that mystery distribution can be. Our goal is to distinguish whether  $X \sim \mu$  or  $X \sim \nu$ . Of course, one can also study the setting in which we receive multiple independent samples  $X_1, \dots, X_n$ .

In this lecture, we consider a special kind of hypothesis testing problem called the *reconstruction problem*, which was devised as basic statistical model for evolutionary genetics; see e.g. [Dur08]. To state the problem, we begin by defining a new type of correlated stochastic process called the *broadcast process*.

**Definition 1** ((Binary Symmetric) Broadcast Process on  $\hat{\mathbb{T}}_d$ ). *Fix  $d \in \mathbb{N}$  and an error parameter  $0 \leq \epsilon \leq 1/2$ . The associated (binary symmetric) broadcast process on the infinite  $d$ -ary tree  $\hat{\mathbb{T}}_d$  rooted at  $r$  is a random assignment  $\sigma_v \in \{\pm 1\}$  to each vertex  $v$  of  $\hat{\mathbb{T}}_d$  generated as follows:*

- *We initialize the process by sampling  $\sigma_r \sim \text{Unif}\{\pm 1\}$  for the root vertex.*
- *Independently for each child  $v$  of  $r$ , we set  $\sigma_v = \sigma_r$  with probability  $1 - \epsilon$  and  $\sigma_v = -\sigma_r$  with probability  $\epsilon$ .*
- *This process continues recursively in the subtrees rooted at each of the children  $v$  of  $r$ .*

*Remark 1.* This process can be vastly generalized to allow for larger state spaces, as well as arbitrary Markov chains for the channels on the edges.

Informally, the vertex  $r$  is “broadcasting” its random assignment to its descendants through channels corrupted by an  $\epsilon$  amount of noise. The basic question is whether or not we can *infer* the assignment of the root vertex given only information about far away vertices. More precisely, if we let  $L(n)$  denote the set of vertices at distance exactly  $n$  from the root, then our goal is to design an estimator  $\hat{\sigma}_{r,n}$  for  $\sigma_r$  based on only the assignments  $\sigma_{L(n)}$  of the vertices in  $L(n)$  for  $n$  large. For instance,  $\hat{\sigma}_{r,n}$  could be a deterministic function mapping  $\{\pm 1\}^{L(n)}$  to  $\{\pm 1\}$  which tells us what to guess upon seeing some specific  $\sigma_{L(n)}$ . At the highest level of generality, our goal is to design a function  $\hat{p} : \{\pm 1\}^{L(n)} \rightarrow [0, 1]$  which outputs the probability that we should set our estimator  $\hat{\sigma}_{r,n} \in \{\pm 1\}$  to, say,  $+1$ .

**Definition 2** (Reconstruction Problem). *Fix  $d \in \mathbb{N}$  and  $0 \leq \epsilon \leq 1/2$ . For an estimator  $\hat{\sigma}_{r,n}$  given access to only  $\sigma_{L(n)}$ , write  $\frac{1+b(\hat{\sigma}_{r,n})}{2} = \Pr[\hat{\sigma}_{r,n} = \sigma_r]$ , where the probability is calculated with respect to a random draw of the broadcast process  $\sigma$  and the randomness of the estimator, and*

$b(\hat{\sigma}_{r,n})$  denotes the “advantage” of the estimator over random guessing.<sup>1</sup> We also write  $\frac{1+b^*}{2}$  for the optimal success probability, where  $b^* = b^*(n, \hat{\mathbb{T}}_d, \epsilon) \geq 0$  is given by maximizing  $b(\hat{\sigma}_{r,n})$  over all possible estimators  $\hat{\sigma}_{r,n}$ .

**Theorem 1.1** (Kesten–Stigum’1966 [KS66], Bleher–Ruiz–Zagreb’1995 [BRZ95]). Fix  $d \in \mathbb{N}$ , and let  $\theta = 1 - 2\epsilon$  for  $0 \leq \epsilon \leq 1/2$  (so that  $\epsilon = \frac{1-\theta}{2}$ ). Let  $\theta_c$  satisfy  $d \cdot \theta_c^2 = 1$ , and correspondingly define  $\epsilon_c = \frac{1}{2} - \frac{1}{2\sqrt{d}}$ . Then we have the following phase transition for the reconstruction problem:

$$\lim_{n \rightarrow \infty} b^*(n, \hat{\mathbb{T}}_d, \epsilon) \begin{cases} = 0, & \text{if } \epsilon > \epsilon_c \text{ (or equivalently } d \cdot \theta^2 < 1) \\ > 0, & \text{if } \epsilon < \epsilon_c \text{ (or equivalently } d \cdot \theta^2 > 1) \end{cases}.$$

*Remark 2.* The reconstruction problem has also been studied on more general classes of trees; see e.g. [Eva+00].

*Remark 3.* It turns out that in the regime  $\epsilon > \epsilon_c$ , there is some  $\delta = \delta(\epsilon) > 0$  such that we have an exponential decay rate

$$b^*(n, \hat{\mathbb{T}}_d, \epsilon) \lesssim (1 - \delta)^n.$$

At criticality, when  $\epsilon = \epsilon_c$ , a closer inspection of the proof reveals that  $b^*(n, \hat{\mathbb{T}}_d, \epsilon) \leq O(1/n)$ . This decay rate is much slower than the decay in the regime  $\epsilon > \epsilon_c$ , but it still yields  $b^*(n, \hat{\mathbb{T}}_d, \epsilon) \rightarrow 0$ .

Very roughly speaking, this phase transition occurs due to two competing effects: There is the exponential-in- $n$  growth of the number of vertices which could receive the “signal”  $\sigma_r$ , but there’s also the exponential decay of information due to the addition of noise in each step as we increase the distance  $n$ .

**A Brief Word on Notation** Throughout the remainder of the course, for a random variable  $X$ , we write  $\text{Law}(X)$  for the associated probability measure  $X$  according to which  $X$  is distributed. We also write  $\text{supp}(X)$  for the *support* of the distribution  $\text{Law}(X)$ , i.e. the set of values which occur with positive probability.

## 2 The Branching Process Perspective

**Theorem 1.1** establishes a phase transition phenomenon for the reconstruction problem on  $\hat{\mathbb{T}}_d$ . Rather than diving head-first into bare-handed calculations, let us first try to relate this to the only phase transition phenomenon we’ve already studied, namely percolation on  $\hat{\mathbb{T}}_d$  (or more generally, Galton–Watson branching processes). Because the broadcast process is about noisy transmission of information, it is natural to track the number of vertices of  $\hat{\mathbb{T}}_d$  which “receive uncorrupted information from the root”. One clean way to formalize this is to reframe the broadcast process as follows:

- We initialize the process by sampling  $\sigma_r \sim \text{Unif}\{\pm 1\}$  for the root vertex.
- Independently for each child  $v$  of  $r$ , we set  $\sigma_v = \sigma_r$  with probability  $\theta = 1 - 2\epsilon$ , and otherwise sample  $\sigma_v \sim \text{Unif}\{\pm 1\}$ .
- This process continues recursively in the subtrees rooted at each of the children  $v$  of  $r$ .

In this way of looking at the broadcast process, we see that independently to each child  $v$ , the parent  $u$  “perfectly transmits” its assignment with probability  $\theta = 1 - 2\epsilon$ . Otherwise, with probability  $2\epsilon$ , the child completely ignores its parent’s assignment and samples a fresh bit in  $\{\pm 1\}$ . In the former case, let us mark the corresponding edge  $\{u, v\}$  with a 1, indicating the edge is “open”; otherwise, we mark the edge  $\{u, v\}$  with a 0, indicating the edge is “closed”. The open edges in  $\hat{\mathbb{T}}_d$ , i.e. the edges which receive a 1, are collectively distributed as bond percolation on  $\hat{\mathbb{T}}_d$  with edge probability  $\theta = 1 - 2\epsilon$ .

<sup>1</sup>A trivial estimator is to just output a uniformly random element of  $\{\pm 1\}$ . This estimator clearly guesses the correct answer with probability  $\frac{1}{2}$ .

For each  $n \in \mathbb{N}$ , let  $Z_n$  denote the number of vertices in  $L(n)$  which are connected to the root  $r$  by a path of open edges (i.e. edges marked by a 1); for each of these vertices, its assignment, as well as the assignments of its ancestors, all perfectly copy the assignment of the root. A key observation is that the collection of random variables  $\mathcal{Z} = \{Z_n\}_{n \in \mathbb{N}}$  is a Galton–Watson branching process with binomial offspring distribution  $\text{Bin}(d, \theta)$ , where recall that  $\theta = 1 - 2\epsilon$ . This branching process is critical if and only if  $d\theta = 1$  and so we expect a phase transition at this point.

In a future lecture, we will discuss the phase transition at  $d\theta = 1$  in the context of the famous (ferromagnetic) *Ising model* from statistical physics, which is intimately related to the broadcast process we’re studying here. However, unfortunately, this transition point doesn’t match the one in [Theorem 1.1](#).

The reason is that even if  $\theta > \frac{1+\epsilon}{d}$ , it could be that the number of vertices  $Z_n$  copying the root is too small relative to the total number of vertices  $d^n$  in  $L(n)$ . In other words, even if there is a “signal” (whose magnitude is  $Z_n$ ) received by  $L(n)$ , it is “drowned out” by the sea of noise from the remaining  $d^n - Z_n$  vertices who receive random bits independent of the root. To get a better sense of “how large” the signal needs to be, imagine an alternative world where all vertices of  $L(n)$  receive independent random bits drawn from  $\text{Unif}\{\pm 1\}$ , and there is simply no way to estimate the root  $\sigma_r$  better than a coin flip. In this “pure-noise” world (sometimes called the *null hypothesis*), the expected number of vertices which have same assignment as the root is  $\frac{1}{2}d^n$  while the standard deviation is  $\frac{1}{2}d^{n/2}$ . In order to “rule out” this alternative universe, we would need  $\mathbb{E}[Z_n] > \frac{1}{2}d^{n/2}$  for large  $n$ . Since  $\mathbb{E}[Z_n] = d^n\theta^n$ , this inequality is equivalent to  $d\theta^2 > 1$ , which is the threshold described by [Theorem 1.1](#). This is a fairly generic phenomenon: The difference in expectations between two distributions needs to exceed the standard deviation of one of them in order to distinguish them. We formalize this in [Section 3.2](#).

### 3 Reinterpreting $b^*$ as a Metric

The problem of recovering  $\sigma_r$  becomes easier if  $\epsilon < \epsilon_c$  is small because the value of  $\sigma_r$  has a noticeable effect on the distribution of  $\sigma_{L(n)}$ . If  $\epsilon < \epsilon_c$  is small and  $\sigma_r = +1$  (resp.  $\sigma_r = -1$ ), then we expect a disproportionate fraction of the vertices in  $L(n)$  to be assigned  $+1$  (resp.  $-1$ ). It turns out we can actually interpret  $b^*$  as measuring the size of this effect. More formally, we can express it as the *total variation distance* between the distribution of  $\sigma_{L(n)}$  conditioned on  $\sigma_r = +1$  and the distribution of  $\sigma_{L(n)}$  conditioned on  $\sigma_r = -1$ . To see this, observe that

$$\begin{aligned} b^* &= 2 \cdot \sup_{\hat{p}} \Pr[\hat{\sigma}_{r,n} = \sigma_r] - 1 \\ &= 2 \cdot \sup_{\hat{p}} \left\{ \frac{1}{2} \Pr[\hat{\sigma}_{r,n} = \sigma_r \mid \sigma_r = +1] + \frac{1}{2} \Pr[\hat{\sigma}_{r,n} = \sigma_r \mid \sigma_r = -1] \right\} - 1 \\ &\quad \text{(Law of Total Probability)} \\ &= \sup_{\hat{p}} \{ \Pr[\hat{\sigma}_{r,n} = +1 \mid \sigma_r = +1] - \Pr[\hat{\sigma}_{r,n} = +1 \mid \sigma_r = -1] \} \\ &= \sup_{\hat{p}} \{ \mathbb{E}_{\sigma} [\hat{p}(\sigma_{L(n)}) \mid \sigma_r = +1] - \mathbb{E}_{\sigma} [\hat{p}(\sigma_{L(n)}) \mid \sigma_r = -1] \} \\ &= \mathcal{D}_{\text{TV}} (\text{Law}(\sigma_{L(n)} \mid \sigma_r = +1), \text{Law}(\sigma_{L(n)} \mid \sigma_r = -1)), \end{aligned} \quad \text{(Lemma 3.1)}$$

where the supremum is over all possible functions  $\hat{p} : \{\pm 1\}^{L(n)} \rightarrow [0, 1]$  and the probability is with respect to  $\sigma$  drawn from the broadcast process and  $\hat{\sigma}_{r,n}$  drawn independently from  $\text{Ber}(\hat{p}(\sigma_{L(n)}))$  given  $\sigma$ . We justify the final step using the following lemma.

**Lemma 3.1** (Total Variation Distance and Test Functions). *Let  $\mu, \nu$  be two probability measures over a common state space  $\Omega$ . Then the total variation distance  $\mathcal{D}_{\text{TV}}(\mu, \nu) \stackrel{\text{def}}{=} \frac{1}{2} \sum_{\omega \in \Omega} |\mu(\omega) - \nu(\omega)|$  may be equivalently written as*

$$\mathcal{D}_{\text{TV}}(\mu, \nu) = \sup_{A \subseteq \Omega} |\mu(A) - \nu(A)| = \sup_f |\mathbb{E}_{\mu}[f] - \mathbb{E}_{\nu}[f]|,$$

where the supremum in the right-most expression is over all 1-bounded functions  $f : \Omega \rightarrow [0, 1]$ .

*Proof.* Since  $\Omega$  is finite, we can view  $f : \Omega \rightarrow [0, 1]$  and  $\mu, \nu$  as big vectors in  $[0, 1]^{\Omega}$ . Then  $|\mathbb{E}_{\mu}[f] - \mathbb{E}_{\nu}[f]| = |\langle \mu - \nu, f \rangle|$ , which is convex in  $f$ . This convexity immediately implies the second equality via the natural correspondence between  $A \subseteq \Omega$  and its  $\{0, 1\}$ -indicator function.

For the first equality, observe that since  $|\mu(A) - \nu(A)| = |\mu(\Omega \setminus A) - \nu(\Omega \setminus A)|$ ,

$$|\mu(A) - \nu(A)| = \frac{1}{2} \left| \sum_{\omega \in A} (\mu(\omega) - \nu(\omega)) \right| + \frac{1}{2} \left| \sum_{\omega \in \Omega \setminus A} (\mu(\omega) - \nu(\omega)) \right| \leq \|\mu - \nu\|_{\text{TV}} \quad \forall A \subseteq \Omega$$

by the Triangle Inequality. On the other hand, this inequality is saturated by the set  $A = \{\omega \in \Omega : \mu(\omega) > \nu(\omega)\}$ , and so we must have equality.  $\square$

### 3.1 Designing Estimators

Given the interpretation of  $b^*$  as the total variation distance between  $\mu_n^+ \stackrel{\text{def}}{=} \text{Law}(\sigma_{L(n)} \mid \sigma_r = +1)$  and  $\mu_n^- \stackrel{\text{def}}{=} \text{Law}(\sigma_{L(n)} \mid \sigma_r = -1)$ , [Lemma 3.1](#) suggests an estimator  $\hat{\sigma}_{r,n}^{\text{MLE}}$  known as the *maximum likelihood estimator*: If  $\sigma_{L(n)} = \tau$  for some  $\tau \in \{\pm 1\}^{L(n)}$ , then

$$\hat{\sigma}_{r,n}^{\text{MLE}} = \begin{cases} +1, & \text{if } \mu_n^+(\tau) > \mu_n^-(\tau) \\ -1, & \text{if } \mu_n^+(\tau) < \mu_n^-(\tau) \end{cases}.$$

This is the optimal estimator whose success probability is precisely  $\frac{1+b^*}{2}$ .

*Remark 4.* Given  $\tau \in \{\pm 1\}^{L(n)}$ , one can of course compute the probabilities  $\mu_n^+(\tau), \mu_n^-(\tau)$  through brute force enumeration of all possible assignments for the vertices in  $L(1), \dots, L(n-1)$ . However, if the goal is to just evaluate whether or not  $\mu_n^+(\tau) > \mu_n^-(\tau)$ , then we can actually do this in time polynomial in the size of the input  $\tau$ . The idea is that by Bayes' Rule,  $\mu_n^+(\tau) > \mu_n^-(\tau)$  holds if and only if

$$\Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau] > \Pr[\sigma_r = -1 \mid \sigma_{L(n)} = \tau].$$

These two quantities can be computed via recursion. The specific algorithm is a *message-passing algorithm* known in the literature as *belief propagation*, the *tree recursion*, or the *sum-product algorithm*, and is extremely well-studied. We will discuss this in greater depth in [Section 5](#); see also [Appendix B](#) for a derivation of this recursion.

Another natural estimator to consider is simply the one given by majority vote:

$$\hat{\sigma}_{r,n}^{\text{MAJ}} \stackrel{\text{def}}{=} \text{sign}(S_n) \quad \text{where} \quad S_n \stackrel{\text{def}}{=} \sum_{v \in L(n)} \sigma_v.$$

We note that the maximum likelihood and majority estimators are not equal.

**Exercise 1.** Prove that  $\hat{\sigma}_{r,n}^{\text{MLE}} \neq \hat{\sigma}_{r,n}^{\text{MAJ}}$ .

While these two estimators have different advantages, i.e.  $b(\hat{\sigma}_{r,n}^{\text{MLE}}) \neq b(\hat{\sigma}_{r,n}^{\text{MAJ}})$ , fortunately for us, it will turn out that  $\lim_{n \rightarrow \infty} b(\hat{\sigma}_{r,n}^{\text{MAJ}}) > 0$  if and only if  $\lim_{n \rightarrow \infty} b(\hat{\sigma}_{r,n}^{\text{MLE}}) > 0$ . This will be a consequence of the proof of [Theorem 1.1](#). We note that there are many other natural broadcast processes for which these two regimes *do not coincide* [[Mos01](#)].

### 3.2 Second Moment Bounds on Total Variation Distance

Before we move to the proof of [Theorem 1.1](#), we will need a technical lemma for bounding the total variation distance. This formalizes the intuition we mentioned in [Section 2](#), where the “signal strength” (i.e. difference in expectations between two hypotheses) needs to exceed the “noise level” (i.e. the typical fluctuations). To state it, for a probability measure  $\mu$  on a state space  $\Omega$  and a function  $f : \Omega \rightarrow \mathbb{R}$ , we abbreviate  $\mathbb{E}_\mu[f] \stackrel{\text{def}}{=} \mathbb{E}_{\omega \sim \mu}[f(\omega)]$  and  $\text{Var}_\mu(f) \stackrel{\text{def}}{=} \text{Var}_{\omega \sim \mu}(f(\omega))$ .

**Lemma 3.2.** *Let  $X, Y$  be two real-valued random variables, and let  $W$  have law  $\frac{1}{2}\text{Law}(X) + \frac{1}{2}\text{Law}(Y)$ .<sup>2</sup> Then we have the lower bound*

$$\mathcal{D}_{\text{TV}}(\text{Law}(X), \text{Law}(Y)) \geq \frac{1}{4} \cdot \frac{(\mathbb{E}[X] - \mathbb{E}[Y])^2}{\text{Var}(W)}.$$

<sup>2</sup>Note that  $W$  can be sampled by first tossing a fair coin, and then setting  $W = X$  (resp.  $W = Y$ ) if the coin comes up heads (resp. tails).

In particular, if  $\mu, \nu$  are two probability measures on a common state space  $\Omega$ , then

$$\mathcal{D}_{\text{TV}}(\mu, \nu) \geq \frac{1}{4} \cdot \sup_{f: \Omega \rightarrow \mathbb{R}} \frac{(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2}{\text{Var}_{\frac{\mu+\nu}{2}}(f)}.$$

For intuition, we also provide a kind of converse to the above.

**Lemma 3.3.** *Let  $\mu, \nu$  be two probability measures on a common state space  $\Omega$ . Then*

$$\mathcal{D}_{\text{TV}}(\mu, \nu)^2 \leq \frac{1}{4} \cdot \sup_{f: \Omega \rightarrow \mathbb{R}} \frac{(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2}{\text{Var}_\mu(f)}.$$

*Remark 5.* For comparison, note that by the Law of Total Variance,

$$\text{Var}_{\frac{\mu+\nu}{2}}(f) = \frac{1}{2} \text{Var}_\mu(f) + \frac{1}{2} \text{Var}_\nu(f) + \frac{1}{4} (\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2.$$

Note that the third term is necessary to ensure the lower bound never exceeds 1, since  $\mathcal{D}_{\text{TV}}(\mu, \nu) \in [0, 1]$  for any  $\mu, \nu$ . So, if  $(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2 \gg \text{Var}_\mu(f), \text{Var}_\nu(f)$ , then Lemma 3.2 yields  $\mathcal{D}_{\text{TV}}(\mu, \nu) \geq 1 - o(1)$ . On the other hand, if  $(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2 \ll \text{Var}_\mu(f), \text{Var}_\nu(f)$ , then Lemma 3.3 yields  $\mathcal{D}_{\text{TV}}(\mu, \nu) \leq o(1)$ .

Proofs for these lemmas are largely based on Cauchy–Schwarz, and provided in Appendix A.

## 4 Proof of Theorem 1.1 in the Case $\epsilon < \epsilon_c$

In this section, we work in the regime  $\epsilon < \epsilon_c$ , or equivalently,  $d\theta^2 > 1$ . From our previous analysis, we have that

$$\begin{aligned} b^* &\geq \mathcal{D}_{\text{TV}}(\text{Law}(\hat{\sigma}_{r,n}^{\text{MAJ}} \mid \sigma_r = +1), \text{Law}(\hat{\sigma}_{r,n}^{\text{MAJ}} \mid \sigma_r = -1)) \\ &= \mathcal{D}_{\text{TV}}(\text{Law}(S_n \mid \sigma_r = +1), \text{Law}(S_n \mid \sigma_r = -1)). \end{aligned}$$

Note the final equality holds due to symmetry between  $\pm 1$ . More concretely, the random variables  $S_n \mid \sigma_r = -1$  and  $-S_n \mid \sigma_r = +1$  have the same law, and so

$$\{s \in \mathbb{Z} : \Pr[S_n = s \mid \sigma_r = +1] > \Pr[S_n = s \mid \sigma_r = -1]\} = \{s \in \mathbb{Z} : s > 0\}$$

and we may apply Lemma 3.1.

We will show that when  $\epsilon < \epsilon_c$ , the right-hand side is lower bounded by a universal constant even in the large  $n$  limit. We will achieve this by combining the second moment method with Lemma 3.2, where  $X_n = S_n \mid \sigma_r = +1$  and  $Y_n = S_n \mid \sigma_r = -1$ ; note that  $W_n = S_n$  without any conditioning. We now compute the quantities involved in Lemma 3.2. If we consider again the Galton–Watson branching process  $\mathcal{Z} = \{Z_n\}_{n \in \mathbb{N}}$  we discussed previously, which has offspring distribution  $\text{Bin}(d, \theta)$  with mean  $d\theta$ , then observe that  $X_n$  is equal to  $Z_n$  plus a sum of  $d^n - Z_n$  many (dependent)  $\text{Unif}\{\pm 1\}$  random variables. Hence,

$$\mathbb{E}[X_n] = \mathbb{E}[Z_n] = d^n \theta^n \quad \text{and} \quad \mathbb{E}[Y_n] = -\mathbb{E}[Z_n] = -d^n \theta^n.$$

We now calculate the variance of  $S_n$ .

$$\begin{aligned} \text{Var}(S_n) &= \mathbb{E}[S_n^2] && \text{(Using } \mathbb{E}[S_n] = 0\text{)} \\ &= \sum_{u, v \in L(n)} \mathbb{E}[\sigma_u \sigma_v] \\ &= \sum_{u \in L(n)} \mathbb{E}[\sigma_u^2] + \sum_{u \in L(n)} \sum_{\ell=0}^{n-1} \sum_{\substack{v \neq u \\ \text{LCA}(u, v) \in L(\ell)}} \mathbb{E}[\sigma_u \sigma_v]. \end{aligned}$$

Fix a pair of distinct vertices  $u, v \in L(n)$  such that their least common ancestor  $w = \text{LCA}(u, v)$  lies in  $L(\ell)$ . Then by the Law of Total Expectation and the fact that an independently sampled  $\text{Unif}\{\pm 1\}$  is assigned when a child fails to copy its parent’s assignment, we have

$$\begin{aligned} \mathbb{E}[\sigma_u \sigma_v] &= \Pr[\sigma_u, \sigma_v \text{ both copy } \sigma_w] \\ &= \Pr[\sigma_u \text{ copies } \sigma_w] \cdot \Pr[\sigma_v \text{ copies } \sigma_w] \\ &= \theta^{2(n-\ell)}. \end{aligned}$$

Combining with the fact that  $\#\{v \neq u : \text{LCA}(u, v) \in L(\ell)\} = d^{n-\ell-1}(d-1)$  for every  $u \in L(n)$ , the above simplifies to

$$\text{Var}(S_n) = d^n + d^{n-1}(d-1) \sum_{k=1}^n (d\theta^2)^k = d^n + d^n \left( (d\theta^2)^n - 1 \right) \cdot \frac{(d-1)\theta^2}{d\theta^2 - 1}.$$

With this in hand, we find that

$$\begin{aligned} b^* &\geq \mathcal{D}_{\text{TV}}(\text{Law}(S_n \mid \sigma_r = +1), \text{Law}(S_n \mid \sigma_r = -1)) && \text{(Previous display)} \\ &\geq \frac{1}{4} \cdot \frac{(\mathbb{E}[X_n] - \mathbb{E}[Y_n])^2}{\text{Var}(S_n)} && \text{(Lemma 3.2)} \\ &= \frac{d^{2n}\theta^{2n}}{d^n + d^n((d\theta^2)^n - 1) \cdot \frac{(d-1)\theta^2}{d\theta^2 - 1}} \\ &= \frac{1}{(d\theta^2)^{-n} + \left(1 - (d\theta^2)^{-n}\right) \cdot \frac{(d-1)\theta^2}{d\theta^2 - 1}} \\ &\rightarrow \frac{d\theta^2 - 1}{(d-1)\theta^2} && \text{(In the } n \rightarrow \infty \text{ limit, using } d\theta^2 > 1\text{)} \\ &= \frac{d}{d-1} \left(1 - (d\theta^2)^{-1}\right). \end{aligned}$$

As a sanity check, this is at most 1 since  $\theta \leq 1$ . It is positive since  $d\theta^2 > 1$  by assumption, which completes the proof of [Theorem 1.1](#).

## 5 Bonus Material: Proof of [Theorem 1.1](#) in the Case $\epsilon > \epsilon_c$

Now suppose  $\epsilon < \epsilon_c$ , i.e.  $d\theta^2 < 1$ . Our goal is to show that  $b^*(n, \widehat{\mathbb{T}}_d, \epsilon) \rightarrow 0$  as  $n \rightarrow \infty$ . Towards this, observe that we can apply Bayes' Rule to obtain that

$$\begin{aligned} &\mathcal{D}_{\text{TV}}(\text{Law}(\sigma_{L(n)} \mid \sigma_r = +1), \text{Law}(\sigma_{L(n)} \mid \sigma_r = -1)) \\ &= \frac{1}{2} \sum_{\tau \in \{\pm 1\}^{L(n)}} |\Pr[\sigma_{L(n)} = \tau \mid \sigma_r = +1] - \Pr[\sigma_{L(n)} = \tau \mid \sigma_r = -1]| \\ &= \frac{1}{2} \sum_{\tau \in \{\pm 1\}^{L(n)}} \left| \Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau] \frac{\Pr[\sigma_{L(n)} = \tau]}{\Pr[\sigma_r = +1]} - \Pr[\sigma_{L(n)} = \tau \mid \sigma_r = -1] \frac{\Pr[\sigma_{L(n)} = \tau]}{\Pr[\sigma_r = -1]} \right| \\ &= \mathbb{E} |\Pr[\sigma_r = +1 \mid \sigma_{L(n)}] - \Pr[\sigma_r = -1 \mid \sigma_{L(n)}]|. \end{aligned}$$

This suggests that we look at the quantity

$$\begin{aligned} \mathcal{M}_n(\tau) &\stackrel{\text{def}}{=} \Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau] - \Pr[\sigma_r = -1 \mid \sigma_{L(n)} = \tau] \\ &= \mathbb{E}[\sigma_r \mid \sigma_{L(n)} = \tau], \quad \forall \tau \in \{\pm 1\}^{L(n)}. \end{aligned}$$

We write  $\mathcal{M}_n$  for the corresponding random variable, where the randomness comes from drawing  $\sigma_{L(n)}$  according to the broadcast process. Due to connections with the (*ferromagnetic*) *Ising model*, this quantity is sometimes referred to as the *magnetization* of the root vertex. Our goal will be to show that  $\mathbb{E}|\mathcal{M}_n|$  decays as  $n \rightarrow \infty$ . Rather than study  $|\mathcal{M}_n|$ , we will instead show step-wise decay for the second moment, which is more tractable.

**Lemma 5.1.** *We have the inequality  $\mathbb{E}|\mathcal{M}_n| \leq \sqrt{\mathbb{E}[\mathcal{M}_n^2]}$ .*

*Proof.* The claim is equivalent to nonnegativity of the variance of  $|\mathcal{M}_n|$ . □

Our goal will be to show that

$$\mathbb{E}[\mathcal{M}_n^2] \leq d\theta^2 \cdot \mathbb{E}[\mathcal{M}_{n-1}^2], \quad (1)$$

since our assumption that  $d\theta^2 < 1$  implies that  $\mathbb{E}|\mathcal{M}_n| \leq \sqrt{\mathbb{E}[\mathcal{M}_n^2]}$  is decaying to zero exponentially fast as  $n \rightarrow \infty$ .

The benefit of switching to the viewpoint of  $\mathcal{M}_n$  is that we can use recursion to understand the probabilities  $\Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau]$  and  $\Pr[\sigma_r = -1 \mid \sigma_{L(n)} = \tau]$ . To state this precisely, let  $u_1, \dots, u_d$  be the children of the root  $r$ , and for each  $i = 1, \dots, d$ , let  $L_i(n-1) \subseteq L(n)$  denote the vertices in the subtree rooted at  $u_i$  which are at distance  $n-1$  from  $u_i$ . For  $\tau \in \{\pm 1\}^{L(n)}$ , we write  $\tau_i \in \{\pm 1\}^{L_i(n-1)}$  for the restriction of  $\tau$  to  $L_i(n-1)$ . Now let  $\mathcal{M}_{n-1}^{(i)}(\tau_i)$  be given by

$$\mathcal{M}_{n-1}^{(i)}(\tau_i) \stackrel{\text{def}}{=} \Pr[\sigma_{u_i} = +1 \mid \sigma_{L_i(n-1)} = \tau_i] - \Pr[\sigma_{u_i} = -1 \mid \sigma_{L_i(n-1)} = \tau_i].$$

Note the corresponding random variables  $\mathcal{M}_{n-1}^{(1)}, \dots, \mathcal{M}_{n-1}^{(d)}$  are correlated through the random root assignment  $\sigma_r$  if  $(\tau_1, \dots, \tau_d) \sim \mu_n$ . However, they become independent if we fix the value of  $\sigma_r$ , i.e. if  $(\tau_1, \dots, \tau_d) \sim \mu_n^{\pm 1}$ . What will be important for us is that  $\text{Law}(\mathcal{M}_{n-1}^{(i)}) = \text{Law}(\mathcal{M}_{n-1})$  individually for each  $i = 1, \dots, d$ . This is because marginally,  $\sigma_{u_i}$  is distributed as  $\text{Unif}\{\pm 1\}$  and the assignment  $\sigma_{L_i(n-1)}$  is independent of  $\sigma_r$  given  $\sigma_{u_i}$ , so  $\text{Law}(\sigma_{L_i(n-1)}) = \text{Law}(\sigma_{L(n-1)})$ .

Finally, let  $f(x) \stackrel{\text{def}}{=} \frac{1-x}{1+x}$ , and for  $d \in \mathbb{N}$ , define

$$g_d(x_1, \dots, x_d) \stackrel{\text{def}}{=} f\left(\prod_{i=1}^d f(x_i)\right).$$

We have the following theorem.

**Theorem 5.2.** *For any  $\tau \in \{\pm 1\}^{L(n)}$ ,*

$$\mathcal{M}_n(\tau) = g_d\left(\theta \cdot \mathcal{M}_{n-1}^{(1)}(\tau_1), \dots, \theta \cdot \mathcal{M}_{n-1}^{(d)}(\tau_d)\right). \quad (2)$$

*In particular, we have the distributional recursion*

$$\mathcal{M}_n \stackrel{d}{=} g_d\left(\theta \cdot \mathcal{M}_{n-1}^{(1)}, \dots, \theta \cdot \mathcal{M}_{n-1}^{(d)}\right). \quad (3)$$

Let us content ourselves for the moment with an informal explanation of how this recursion is derived; we formally prove this in [Appendix B](#). The idea is that because  $\sigma_{L_i(n-1)}$  is independent of  $\sigma_r$  given the value of  $\sigma_{u_i}$ , we can express  $\Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau]$  as a function of the probabilities  $\Pr[\sigma_{u_i} = \pm 1 \mid \sigma_{L_i(n-1)} = \tau_i]$ , for any  $\tau \in \{\pm 1\}^{L(n)}$ . The specific function is known in the literature as *belief propagation*, the *tree recursion*, or the *sum-product algorithm*, and is extremely well-studied. Typically, it is discussed in the context of spin systems like the Ising model; we derive it using Bayes' Rule in [Appendix B](#) and postpone a more in-depth treatment of spin systems to a future lecture.

Taking [Theorem 5.2](#) as a given, to establish an inequality like [Eq. \(1\)](#), our goal is now to bound the rather complicated right-hand side of [Eq. \(3\)](#) by something like

$$\mathbb{E}\left[f\left(\prod_{i=1}^d f\left(\theta \cdot \mathcal{M}_{n-1}^{(i)}\right)\right)\right]^2 \leq \theta^2 \sum_{i=1}^d \mathbb{E}\left[\left(\mathcal{M}_{n-1}^{(i)}\right)^2\right] = d\theta^2 \cdot \mathbb{E}[\mathcal{M}_{n-1}^2]. \quad (4)$$

With this as our aim, let us begin by massaging [Eq. \(3\)](#). We build up  $g_d$  inductively.

**Lemma 5.3.** *For every  $d \in \mathbb{N}$ , we have that*

$$g_d(x_1, \dots, x_d) = g_2(g_{d-1}(x_1, \dots, x_{d-1}), x_d),$$

*where  $g_2$  is explicitly given by the expression*

$$g_2(x, y) = \frac{1 - \frac{1-x}{1+x} \cdot \frac{1-y}{1+y}}{1 + \frac{1-x}{1+x} \cdot \frac{1-y}{1+y}} = \frac{x+y}{1+xy}.$$

*Proof.* This is a straightforward calculation, taking advantage of the fact that  $f(f(x)) = x$ .  $\square$

The combinatorial interpretation of the conjunction of [Theorem 5.2](#) and [Lemma 5.3](#) is the following. Imagine we inductively build up our infinite  $d$ -ary tree  $\widehat{\mathbb{T}}_d$  rooted at  $r$  by first building  $d$  trees  $T_1, \dots, T_d$ , where each tree  $T_i$  is an infinite  $d$ -ary tree rooted at  $u_i$  along with a single edge



joining  $u_i$  to a special leaf vertex  $r_i$ . Independently in each  $T_i$ , we can run the same broadcast process initialized from  $r_i$ . A direct calculation reveals that  $\theta \cdot \mathbb{E} \left[ \mathcal{M}_{n-1}^{(i)} \right]$  is precisely the total variation distance/optimal advantage  $b^*(n, T_i, \epsilon)$  with respect to  $T_i$ . If we merge the vertices  $r_1, \dots, r_d$  into a single vertex, then we recover  $\mathbb{T}_d$ . The purpose of the function  $g_d$  is to compute the effect this merge operation induces on the reconstruction probability for the root  $r$ . [Lemma 5.3](#) captures the step-by-step effect of merging  $r_1, r_2$ , then merging  $r_1, r_2, r_3$ , etc.

Towards an inequality analogous to [Eq. \(4\)](#), we would like to establish that the factor  $\frac{1}{1+xy}$  is at most 1 in absolute value, at least in expectation when we plug in our random variables. This expression looks unwieldy, so let us attempt to “linearize”  $g_2$ . Using the identity  $\frac{1}{1+r} = 1 - r + \frac{r^2}{1+r}$ , we have that

$$\begin{aligned} g_2(x, y) &= (x + y) \left( 1 - xy + \frac{x^2 y^2}{1 + xy} \right) = (x + y) - x^2 y - xy^2 + x^2 y^2 \cdot g_2(x, y) \\ &\leq (x + y) - x^2 y - xy^2 + x^2 y^2, \end{aligned} \quad (\text{For } x, y \geq -1)$$

where in the final step, we used the fact that  $f([-1, +\infty]) = \mathbb{R}_{\geq 0} \cup \{+\infty\}$ ; in particular,  $-1$  is not in the range of  $f$  in our context. Our goal will be to ensure that the terms  $-x^2 y - xy^2 + x^2 y^2$  are negative, at least in expectation, when we plug in our random variables. For this, we need one final rather curious set of identities.

**Lemma 5.4.** *Recalling the notation  $\mu_n^+ = \text{Law}(\sigma_{L(n)} \mid \sigma_r = +1)$  and  $\mu_n = \text{Law}(\sigma_{L(n)})$ , we have the identities*

$$\begin{aligned} \mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} [\mathcal{M}_n] &= \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} [\mathcal{M}_n^2] \\ \mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} [\mathcal{M}_{n-1}^{(i)}] &= \theta \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2] \\ \mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(i)} \right)^2 \right] &= \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2], \end{aligned}$$

for any  $i = 1, \dots, d$ .

We prove [Lemma 5.4](#) at the end of this section. With these tools in hand, we prove our final contraction inequality. We do the case  $d = 2$  here (with suppressed notation); the general case follows by a straightforward induction, using [Lemma 5.3](#) to add one child subtree at a time. We have

$$\begin{aligned} \mathbb{E}_{\mu_n} [\mathcal{M}_n^2] &= \mathbb{E}_{\mu_n^+} [\mathcal{M}_n] && (\text{Lemma 5.4}) \\ &= \mathbb{E}_{\mu_n^+} \left[ g_2 \left( \theta \cdot \mathcal{M}_{n-1}^{(1)}, \theta \cdot \mathcal{M}_{n-1}^{(2)} \right) \right] && (\text{Using Eq. (2)}) \\ &\leq \theta \cdot \mathbb{E}_{\mu_n^+} [\mathcal{M}_{n-1}^{(1)}] + \theta \cdot \mathbb{E}_{\mu_n^+} [\mathcal{M}_{n-1}^{(2)}] \\ &\quad - \theta^3 \cdot \mathbb{E}_{\mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(1)} \right)^2 \right] \cdot \mathbb{E}_{\mu_n^+} [\mathcal{M}_{n-1}^{(2)}] - \theta^3 \cdot \mathbb{E}_{\mu_n^+} [\mathcal{M}_{n-1}^{(1)}] \cdot \mathbb{E}_{\mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(2)} \right)^2 \right] \\ &\quad + \theta^4 \cdot \mathbb{E}_{\mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(1)} \right)^2 \right] \cdot \mathbb{E}_{\mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(2)} \right)^2 \right] && (\text{Independence under } \mu_n^+) \\ &= 2\theta^2 \cdot \mathbb{E}_{\mu_{n-1}} [\mathcal{M}_{n-1}^2] - \theta^4 \cdot \mathbb{E}_{\mu_{n-1}} [\mathcal{M}_{n-1}^2]^2 && (\text{Lemma 5.4}) \\ &\leq 2\theta^2 \cdot \mathbb{E}_{\mu_{n-1}} [\mathcal{M}_{n-1}^2]. \end{aligned}$$

All that remains to complete the proof of [Theorem 1.1](#) in the  $\epsilon > \epsilon_c$  regime is to prove [Lemma 5.4](#) and [Theorem 5.2](#). The former is proved here, while the latter is relegated to [Appendix B](#).

*Proof of Lemma 5.4.* For this first identity, observe that

$$\begin{aligned} \mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} [\mathcal{M}_n] &= \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} \left[ \frac{\Pr[\sigma_{L(n)} \mid \sigma_r = +1]}{\Pr[\sigma_{L(n)}]} (\Pr[\sigma_{L(n)} \mid \sigma_r = +1] - \Pr[\sigma_{L(n)} \mid \sigma_r = -1]) \right] \\ &\quad (\text{Change of densities}) \\ &= \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} [2 \Pr[\sigma_r = +1 \mid \sigma_{L(n)}] (\Pr[\sigma_{L(n)} \mid \sigma_r = +1] - \Pr[\sigma_{L(n)} \mid \sigma_r = -1])] \\ &\quad (\text{Bayes' Rule}) \\ &= \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} [\mathcal{M}_n^2 + \mathcal{M}_n] && (*) \\ &= \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} [\mathcal{M}_n^2]. && (\text{Using } \mathbb{E}_{\sigma_{L(n)} \sim \mu_n} [\mathcal{M}_n] = 0 \text{ by symmetry}) \end{aligned}$$



For  $(*)$ , we used the fact that  $\Pr[\sigma_r = +1 \mid \sigma_{L(n)}] + \Pr[\sigma_r = -1 \mid \sigma_{L(n)}] = 1$  implies

$$2\Pr[\sigma_r = +1 \mid \sigma_{L(n)}] = 1 + (\Pr[\sigma_r = +1 \mid \sigma_{L(n)}] - \Pr[\sigma_r = -1 \mid \sigma_{L(n)}]).$$

This proves the first identity. For the second identity, we use the Law of Total Expectation to obtain

$$\begin{aligned}\mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} [\mathcal{M}_{n-1}^{(i)}] &= \theta \cdot \mathbb{E} [\mathcal{M}_{n-1}^{(i)} \mid u_i \text{ copies } r \text{ directly}] + (1 - \theta) \cdot \mathbb{E} [\mathcal{M}_{n-1}^{(i)} \mid \sigma_{u_i} \sim \text{Unif}\{\pm 1\}] \\ &= \theta \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}^+} [\mathcal{M}_{n-1}] + (1 - \theta) \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}] \\ &= \theta \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2].\end{aligned}\quad (\text{First identity})$$

For the final identity, we again use the Law of Total Expectation, similar to the above, to obtain

$$\mathbb{E}_{\sigma_{L(n)} \sim \mu_n^+} \left[ \left( \mathcal{M}_{n-1}^{(i)} \right)^2 \right] = \theta \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}^+} [\mathcal{M}_{n-1}^2] + (1 - \theta) \cdot \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2].$$

But  $\mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}^+} [\mathcal{M}_{n-1}^2] = \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}^-} [\mathcal{M}_{n-1}^2] = \mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2]$  by symmetry of the assignments  $\pm 1$ , and so the above is equal to  $\mathbb{E}_{\sigma_{L(n-1)} \sim \mu_{n-1}} [\mathcal{M}_{n-1}^2]$ .  $\square$

## References

- [Bor+06] Christian Borgs, Jennifer Chayes, Elchanan Mossel, and Sebastien Roch. “The Kesten-Stigum Reconstruction Bound Is Tight for Roughly Symmetric Binary Channels”. In: *Proceedings of the 47th Annual IEEE Symposium on Foundations of Computer Science. FOCS '06*. USA: IEEE Computer Society, 2006, pp. 518–530. ISBN: 0769527205. DOI: [10.1109/FOCS.2006.76](#) (cit. on p. 1).
- [BRZ95] P. M. Bleher, J. Ruiz, and V. A. Zagrebnov. “On the purity of the limiting gibbs state for the Ising model on the Bethe lattice”. In: *Journal of Statistical Physics* 79 (1995), pp. 473–482 (cit. on p. 2).
- [Dur08] Richard Durrett. *Probability Models for DNA Sequence Evolution*. Probability and Its Applications. Springer New York, NY, 2008 (cit. on p. 1).
- [Eva+00] William Evans, Claire Kenyon, Yuval Peres, and Leonard J. Schulman. “Broadcasting on trees and the Ising model”. In: *The Annals of Applied Probability* 10.2 (2000), pp. 410–433 (cit. on p. 2).
- [KS66] H. Kesten and B. P. Stigum. “Additional Limit Theorems for Indecomposable Multi-dimensional Galton-Watson Processes”. In: *The Annals of Mathematical Statistics* 37.6 (1966), pp. 1463–1481 (cit. on p. 2).
- [Mos01] Elchanan Mossel. “Reconstruction on Trees: Beating the Second Eigenvalue”. In: *The Annals of Applied Probability* 11.1 (2001), pp. 285–300 (cit. on p. 4).

## A Proofs for Second Moment Total Variation Bounds

*Proof of Lemma 3.2.* Without loss of generality, we may shift  $X, Y, W$  by the same constant to ensure they have mean zero, while preserving the total variation distance between  $\text{Law}(X), \text{Law}(Y)$ . Hence, we assume  $\mathbb{E}[X] = \mathbb{E}[Y] = \mathbb{E}[W] = 0$  in the remainder of the proof.

Observe that  $\text{supp}(W) = \text{supp}(X) \cup \text{supp}(Y)$ , and so

$$\begin{aligned}\mathcal{D}_{\text{TV}}(\text{Law}(X), \text{Law}(Y)) &= \frac{1}{2} \sum_{\omega \in \text{supp}(W)} |\Pr[X = \omega] - \Pr[Y = \omega]| \\ &= \sum_{\omega \in \text{supp}(W)} \frac{|\Pr[X = \omega] - \Pr[Y = \omega]|}{2 \cdot \Pr[W = \omega]} \cdot \Pr[W = \omega] \\ &= \mathbb{E}[|f(W)|],\end{aligned}$$

where  $f(\omega) \stackrel{\text{def}}{=} \frac{\Pr[X=\omega] - \Pr[Y=\omega]}{\Pr[X=\omega] + \Pr[Y=\omega]}$  takes values in the interval  $[-1, 1]$ . Hence, we may lower bound the above by

$$\begin{aligned} \mathbb{E}[f(W)^2] &\geq \frac{\mathbb{E}[W \cdot f(W)]^2}{\mathbb{E}[W^2]} && \text{(Cauchy-Schwarz)} \\ &= \frac{1}{\text{Var}(W)} \cdot \left( \sum_{\omega \in \text{supp}(W)} \omega \cdot \Pr[W = \omega] \cdot \frac{\Pr[X = \omega] - \Pr[Y = \omega]}{2 \cdot \Pr[W = \omega]} \right)^2 \\ &&& (\mathbb{E}[W] = 0 \text{ and definition of } f) \\ &= \frac{1}{4} \cdot \frac{(\mathbb{E}[X] - \mathbb{E}[Y])^2}{\text{Var}(W)}. \end{aligned}$$

For the second claim, observe that  $f : \Omega \rightarrow \mathbb{R}$  is any function, then letting  $X = f(x)$  for  $x \sim \mu$  and  $Y = f(y)$  for  $y \sim \nu$ , we have

$$\mathcal{D}_{\text{TV}}(\mu, \nu) \geq \mathcal{D}_{\text{TV}}(\text{Law}(X), \text{Law}(Y)) \geq \frac{1}{4} \cdot \frac{(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2}{\text{Var}_{\frac{\mu+\nu}{2}}(f)}.$$

Since  $f$  was arbitrary, the proof is complete.  $\square$

*Proof of Lemma 3.3.* The key idea is to use an intermediate quantity, namely the  $\chi^2$ -squared divergence, which is defined as

$$\chi^2(\nu \parallel \mu) \stackrel{\text{def}}{=} \mathbb{E}_{\omega \sim \mu} \left[ \left( 1 - \frac{\nu(\omega)}{\mu(\omega)} \right)^2 \right].$$

Note that this is simply the variance of the function  $\frac{d\nu}{d\mu}(\omega) \stackrel{\text{def}}{=} \frac{\nu(\omega)}{\mu(\omega)}$ , often referred to as the *density* of  $\nu$  with respect to  $\mu$ . With this, we have by Cauchy-Schwarz that

$$\mathcal{D}_{\text{TV}}(\mu, \nu)^2 = \frac{1}{4} \cdot \mathbb{E}_{\omega \sim \mu} \left[ \left| 1 - \frac{\nu(\omega)}{\mu(\omega)} \right| \right]^2 \leq \frac{1}{4} \cdot \chi^2(\nu \parallel \mu).$$

On the other hand, we claim that the  $\chi^2$ -squared divergence admits the variational formula

$$\chi^2(\nu \parallel \mu) = \sup_{f: \Omega \rightarrow \mathbb{R}} \frac{(\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2}{\text{Var}_\mu(f)}.$$

To see this, observe that for any  $f$ , we have

$$\begin{aligned} \mathbb{E}_\mu[f] - \mathbb{E}_\nu[f] &= \mathbb{E}_\nu(f - \mathbb{E}_\mu[f]) \\ &= \mathbb{E}_\mu \left[ \frac{d\nu}{d\mu} \cdot (f - \mathbb{E}_\mu[f]) \right] && \text{(Change of Measure)} \\ &= \mathbb{E}_\mu \left[ \left( \frac{d\nu}{d\mu} - 1 \right) \cdot (f - \mathbb{E}_\mu[f]) \right] && \text{(Using } \mathbb{E}_\mu[f - \mathbb{E}_\mu[f]] = 0) \end{aligned}$$

It follows that

$$\begin{aligned} (\mathbb{E}_\mu[f] - \mathbb{E}_\nu[f])^2 &= \mathbb{E}_\mu \left[ \left( \frac{d\nu}{d\mu} - 1 \right) \cdot (f - \mathbb{E}_\mu[f]) \right]^2 \\ &\leq \mathbb{E}_\mu \left[ \left( \frac{d\nu}{d\mu} - 1 \right)^2 \right] \cdot \mathbb{E}_\mu [(f - \mathbb{E}_\mu[f])^2] && \text{(Cauchy-Schwarz)} \\ &= \text{Var}_\mu \left( \frac{d\nu}{d\mu} \right) \cdot \text{Var}_\mu(f) \\ &= \chi^2(\nu \parallel \mu) \cdot \text{Var}_\mu(f). && \text{(Definition of } \chi^2(\nu \parallel \mu)) \end{aligned}$$

Rearranging proves that  $\chi^2(\nu \parallel \mu)$  is lower bounded by the given variational formula. To prove equality, we need to exhibit a function  $f : \Omega \rightarrow \mathbb{R}$  which achieves equality. For this, we simply use the density  $\frac{d\nu}{d\mu}$ . Plugging in this function, we have

$$\mathbb{E}_\mu \left[ \frac{d\nu}{d\mu} \right] - \mathbb{E}_\nu \left[ \frac{d\nu}{d\mu} \right] = 1 - \mathbb{E}_\mu \left[ \left( \frac{d\nu}{d\mu} \right)^2 \right] = -\text{Var}_\mu \left( \frac{d\nu}{d\mu} \right) = -\chi^2(\nu \parallel \mu).$$

Plugging this back into the ratio in the variational expression completes the proof.  $\square$

## B Proof of Theorem 5.2

We express  $\Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau]$  in terms of the quantities  $\Pr[\sigma_{u_i} = \pm 1 \mid \sigma_{L_i(n-1)} = \tau_i]$ . To do this, first observe that because the assignments  $\sigma_{L_1(n-1)}, \dots, \sigma_{L_d(n-1)}$  are independent conditioned on the value of  $\sigma_r$ , we have

$$\Pr[\sigma_{L(n)} = \tau \mid \sigma_r = +1] = \prod_{i=1}^d \Pr[\sigma_{L_i(n-1)} = \tau_i \mid \sigma_r = +1].$$

Now, using the Law of Total Probability,

$$\begin{aligned} & \Pr[\sigma_{L_i(n-1)} = \tau_i \mid \sigma_r = +1] \\ &= (1 - \epsilon) \cdot \Pr[\sigma_{L_i(n-1)} = \tau_i \mid \sigma_{u_i} = +1] + \epsilon \cdot \Pr[\sigma_{L_i(n-1)} = \tau_i \mid \sigma_{u_i} = -1] \\ &= 2 \cdot \Pr[\sigma_{L_i(n-1)} = \tau_i] \cdot \left( 1 + \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right). \end{aligned}$$

These calculations all hold mutatis mutandis for the case  $\sigma_r = -1$ . Since

$$\Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau] = \frac{\Pr[\sigma_{L(n)} = \tau \mid \sigma_r = +1]}{\Pr[\sigma_{L(n)} = \tau \mid \sigma_r = +1] + \Pr[\sigma_{L(n)} = \tau \mid \sigma_r = -1]},$$

which follows by Bayes' Rule and  $\Pr[\sigma_r = +1] = \Pr[\sigma_r = -1] = 1/2$ , plugging in the above calculations and canceling out the factors of  $2 \cdot \Pr[\sigma_{L_i(n-1)} = \tau_i]$  yield

$$\begin{aligned} \Pr[\sigma_r = +1 \mid \sigma_{L(n)} = \tau] &= \frac{\prod_{i=1}^d \left( 1 + \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right)}{\prod_{i=1}^d \left( 1 + \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right) + \prod_{i=1}^d \left( 1 - \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right)} \\ \Pr[\sigma_r = -1 \mid \sigma_{L(n)} = \tau] &= \frac{\prod_{i=1}^d \left( 1 - \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right)}{\prod_{i=1}^d \left( 1 + \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right) + \prod_{i=1}^d \left( 1 - \theta \cdot \mathcal{M}_{n-1}^{(i)}(\tau_i) \right)}. \end{aligned}$$

This is one way of expressing the *belief propagation* equations, although they are more commonly written in terms of  $\epsilon$  and the conditional marginal probabilities  $\Pr[\sigma_{u_i} = \pm 1 \mid \sigma_{L_i(n-1)} = \tau_i]$  rather than  $\theta$  and  $\mathcal{M}_{n-1}^{(i)}(\tau_i)$ . Subtracting these two expressions and dividing both the numerator and denominator by common factors then yield the recursion for  $\mathcal{M}_n(\tau)$  stated in Theorem 5.2.