





Article

Bayesian Network Analysis of Lysine Biosynthesis Pathway in Rice

Aditya Lahiri ^{1,†}, Khushboo Rastogi ², Aniruddha Datta ¹ and Endang M. Septiningsih ²

¹ Department of Electrical and Computer Engineering, Texas A&M University, 188 Bizzell St, College Station, TX 77843, U.S.A.; adi441994@gmail.com (A.L.); datta@ece.tamu.edu (A.D.)

² Department of Soil and Crop Sciences, Texas A&M University, College Station, TX 77843, U.S.A.; khushboorastogi5@tamu.edu (K.R.); eseptiningsih@tamu.edu (E.S.)

* Correspondence: adi441994@gmail.com; Tel.: +1-832-782-4580

Abstract: Lysine is the first limiting essential amino acid in rice because it is present in the lowest quantity compared to all the other amino acids. Amino acids are the building block of proteins and play an essential role in maintaining the human body's healthy functioning. Rice is a staple food for large proportion of the global population, thus increasing the lysine content in rice will improve its nutritional value. In this paper, we studied the lysine biosynthesis pathway in rice (*Oryza Sativa*) to identify the regulators of the lysine reporter gene *LYSA* (LOC_Os02g24354). Genetically intervening at the regulators has the potential to increase the overall lysine content in rice. We modeled the lysine biosynthesis pathway in rice seedlings under normal and saline (NaCl) stress conditions using Bayesian networks. We estimated the model parameters using experimental data and identified the gene *DAPF* (LOC_Os12g37960) as a positive regulator of the lysine reporter gene *LYSA* under both normal and saline stress conditions. Based on this analysis, we conclude that the gene *DAPF* is a potent candidate for genetic intervention. Upregulating *DAPF* using methods such as CRISPR-Cas9 has the potential to upregulate the lysine reporter gene *LYSA* and increase the overall lysine content in rice.

Keywords: Lysine; Rice; Amino Acids; Saline Stress; Abiotic Stress; Gene Regulatory Network; Bayesian Network; Parameter Estimation; Inference; RNA Seq

1. Introduction

1.1. Background

Proteins are one of the primary building blocks of all life on Earth and are present in every cell in the human body. Proteins are a crucial macro-nutrient in the human diet; they help build and repair cells and are essential for the human body's growth and development [1]. Proteins are comprised of long chains of amino acids; once the human body digests the proteins, they are broken down into their constituent amino acids [2]. There are twenty naturally existing amino acids that encode the 20,000 (approximate) unique proteins in the human body [3]. Among these amino acids, nine are classified as essential, and eleven are classified as non-essential [2,3]. Amino acids produced by the human body are considered non-essential, whereas the amino acids that cannot be synthesized by the body are considered essential [3]. Essential amino acids include phenylalanine, valine, tryptophan, threonine, isoleucine, methionine, histidine, leucine, and lysine [4]. Since essential amino acids cannot be synthesized, they need to be introduced to the human body through diets rich in *complete proteins*. A protein food source is considered a complete protein if it contains all the essential amino acids [5]. Typically animal-based proteins are considered sources of complete protein. On the other hand, plant-based proteins are considered incomplete as they do not contain all the essential amino acids [5,6].

According to the National Academy of Medicine, the recommended dietary allowance (RDA) of protein intake is 0.8 g/kg/day [7,8]. A diet deficient in protein can cause edema, thinning of hair, and muscle mass loss in adults [9]. Though protein deficiency is rare in the developed world, it is still prevalent in impoverished and underdeveloped countries, especially among children [9,10]. Plant-based proteins accounted for 57 % of the global

protein supply and were followed by protein sourced meat and dairy, which accounted for 18% and 10%, respectively [11]. Even though plant-based proteins constitute a majority of the global protein supply, according to the World Health Organization (WHO) the demand for animal-based protein has been on the rise due to urbanization, population growth, and rising economies. The WHO predicts the annual meat production to reach 376 million tonnes by 2030, a 72% increase since 1997-1999 when the yearly meat production was 218 million tonnes [12]. This global increase has placed a burden on the livestock sector, especially in Europe and the Americas, where animal-based protein intake is higher than plant-based proteins [13]. In the USA and European countries, proteins from animal-based sources ranged from 55% to 71 % (depending on countries) of the total protein intake, a significant proportion of which were from red meat [14].

Animal-based protein sources such as meat, milk, and eggs are richer in essential amino acids and have a higher food protein quality in terms of digestibility, net protein utilization, and biological value compared to plant-based protein sources like legumes and cereals [13]. However, animal-based proteins, specifically processed and red meats, have been linked with cancer, type 2 diabetes, and cardiovascular diseases [15–17]. Apart from health concerns, proteins sourced from animals have a significant impact on climate change. According to the Food and Agriculture Organization of the United Nations, the livestock supply chain accounts for 14.5% of global anthropogenic greenhouse gas emissions [18]. With the global population set to reach 9.8 billion by 2050 and the increasing demand for animal-based proteins, the challenges associated with food security and climate change will only be exacerbated [12,19]. Hence, a shift towards plant-based protein sources may help reduce the carbon footprint, risks of chronic illness, and food security. While plant-based proteins may not contain all the necessary essential amino acids, a diet containing a diverse range of plant proteins can help overcome this limitation [20]. Cereal plants such as wheat, rice, and maize constitute the primary protein sources in developing countries [21,22]. With the majority of the world's population living in developing countries, it will therefore be beneficial to increase the protein content in cereal plants to ensure food security and prevent malnutrition.

1.2. Lysine Content in Rice

Lysine is produced in the aspartate pathway along with three other essential amino acids threonine, methionine, and isoleucine [23]. Lysine is also the first limiting essential amino acid in cereal and legume crops because it is present in the lowest quantity [23–25]. This is why lysine deficiency is a common problem in developing nations that rely heavily on cereal crops [23,26]. A lysine deficient diet can reduce immunity, decrease protein levels in the blood, and cause retardation of mental and physical development in children [24]. Rice is a cereal plant that is an important food source for more than 50% of the global population [27]. About 95% of global rice is produced in developing countries, among which 92% are countries in Asia [28]. Rice accounts for 50% of the dietary caloric supply for 520 million living in poverty in Asia [29]. Like most cereal crops, rice is deficient in lysine, so in this study, we are interested in identifying the genetic regulators of lysine in rice, since intervening at these regulators has the potential to increase the free lysine content in rice [30]. Enriching lysine content in rice will be a step towards ensuring food security and preventing malnutrition especially in the vulnerable sectors of the global population.

Over the last 50 years, lysine metabolism has been extensively studied. It has been shown that lysine is a self-regulating amino acid as the lysine biosynthesis pathway has two inhibition feedback loops. These feedback loops are activated by the free lysine content, which negatively regulates the enzymes dihydrodipicolinate synthase (DHPS) and aspartate kinase (AK) [24,31]. AK is the first enzyme of the lysine biosynthesis pathway and is also inhibited by threonine, another essential amino acid synthesized by the aspartate pathway [24,31]. Lysine is also degraded through the enzymes lysine ketoglutarate reductase (LKR) and saccharopine dehydrogenase (SDH) bifunctional enzymes [31]. The LKR and SDH enzymes are present in the saccharopine pathway and they initiate the

lysine catabolism process through the TCA cycle (tricarboxylic acid cycle) [24]. Thus lysine can be enriched in cereal plants by enhancing its production in the biosynthesis pathway, preventing its catabolism, or combining these two approaches. A study by Long *et al.* (2013) focused on enhancing lysine through metabolic engineering of rice. These transgenic lines of rice overexpressed AK and DHPS. They observed that LKR and SDH levels were significantly higher in seeds of these rice lines, implying that the catabolic enzymes LKR and SDH were counteracting the effects of transgene AK and DHPS [32]. This method increased the free lysine content by 1.1 times in transgenic lines compared to the wild type. This study also implemented a LKR-RNAi line, which showed a 10 fold increase in lysine content, and a combination LKR-RNAi with AK/DHPS overexpressing lines led to a 60 fold increase in free lysine content. In a different study, Yang *et al.* (2016) developed two pyramid transgenic lines in rice. The lysine content in these transgenic lines showed increased lysine content up to 25 fold. This was achieved by enhancing the biosynthesis pathway and suppressing the catabolism pathway at the same time [33]. Unlike many lysine enhancement studies, which lead to reduced yield, oil content, and phenotype change, no significant trait changes were observed in this case, and the developed transgenic rice was deemed favorable for commercialization [34–36].

While these studies have demonstrated that lysine content can be enhanced through careful metabolic engineering of high-lysine transgenic lines, these are not yet commercialized. Furthermore, transgenic crops rely on introducing foreign genes (transgenes) into the host crop, making them vulnerable to public acceptance. That is why in this paper, we are interested in understanding the underlying genetic regulatory networks (GRNs) that govern these complex interactions. The GRNs can help us identify the genetic regulators of lysine which can be targeted using gene-editing methods such as CRISPR-Cas9. Unlike transgenic crops, the final product of gene editing can be cleared of any foreign DNA segments. Instead of relying on transgenic insertions, gene editing instead knocks out or replaces targeted native genes in the genome of the crop to give rise to desirable traits. The United States Department of Agriculture (USDA) has allowed gene edited crops to be labeled as non-GMO, which will make gene edited crops significantly less controversial than transgenic crops [37]. A recent study by Shew *et al.* showed that gene edited crops were preferred over GMO crops in multiple countries [38]. Thus by studying the underlying GRN involved in lysine regulation in rice we can identify potential targets for gene editing.

LKR and SDH are known regulators of lysine in the catabolic pathway, and genetically intervening at them can prevent lysine degradation [39]. Therefore in this paper, we focus on identifying lysine regulators in the biosynthesis pathway. Overexpressing the regulators in the biosynthesis pathway through gene editing techniques such as CRISPR-Cas9 has the potential to increase the free lysine content in rice. In Figure 1, we derive the GRN of the lysine biosynthesis pathway in rice (*Oryza Sativa*) from the KEGG pathways database [40]. Each rectangular box in Figure 1 represents a gene in the lysine biosynthesis pathway. The gene names are annotated according to their respective MSU IDs (LOC_Os###g#####) [41]. In addition to the MSU IDs, the boxes contain alphabets in red font within parenthesis. These alphabets are used as an alias for genes in the later sections of the paper. Genes I-N have been given names in the literature and these names have been mentioned in the boxes alongside their MSU IDs, for e.g. gene K (LOC_Os03g09910) is also known as *ALD1*. The genetic interactions converge at *LYSA* (LOC_Os02g24354 or gene N) which positively regulates the amino acid lysine (L-Lysine, where the α carbon is in the S configuration). This makes *LYSA* (gene N) a reporter gene of lysine. Thus our objective is to identify genes that will upregulate *LYSA*.

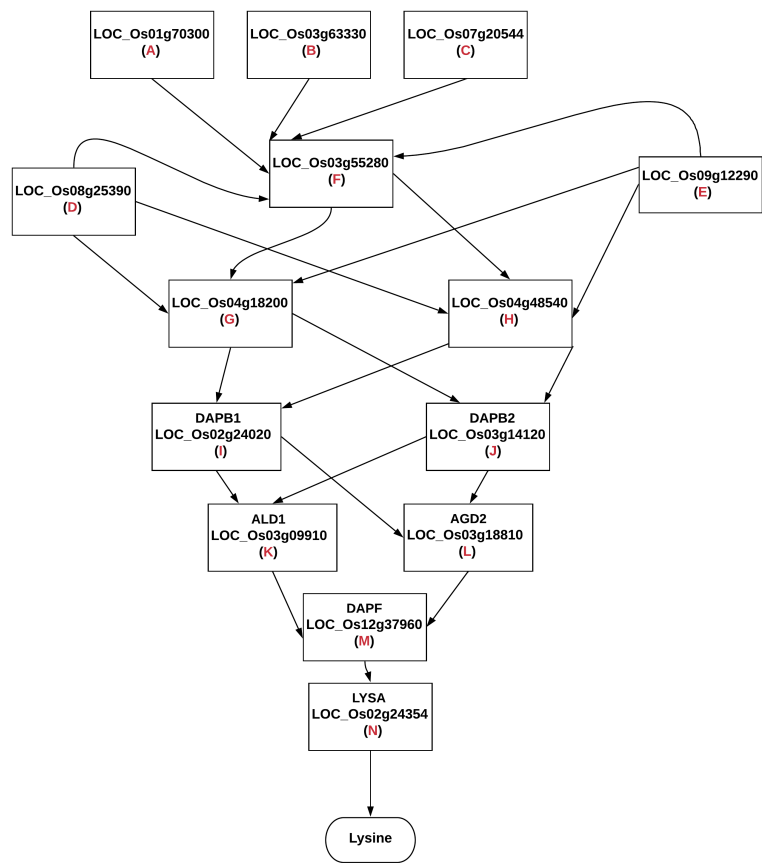


Figure 1. Gene regulatory network for lysine biosynthesis pathway in rice. The gene names are presented according to their MSU IDs. The alphabets in red font are aliases for the respective genes. For e.g. LOC_Os01g70300 will be referred to as gene A. Genes I-N have been given names in the literature, these have been mentioned in the figure alongside their respective MSU IDs.

To identify the *LYSA* regulators, we will model the GRN of the lysine biosynthesis pathway using Bayesian Networks (BN). We will then use publicly available data to infer the BN model’s parameters. The model can then be used to identify the genes that upregulate *LYSA*. This modeling pipeline is similar to our previous work where we identified regulators of drought response in *Arabidopsis* [42,43]. We identify the *LYSA* regulators under normal and saline stress (NaCl) conditions. Soil salinity is one of the significant environmental constraints on the crop life cycle. Nearly 5% (77 million hectares) of the global arable land has excess salinity [44]. Due to various factors such as climate change and irrigation malpractices, the soil salinity is predicted to increase by 16.2 million hectares by 2050 [45,46]. Among abiotic stresses, soil salinity is the second largest cause of crop loss in rice after drought [47,48]. Saline stress primarily affects rice during its seedling, early vegetative, and reproductive stages [47,49]. We have extensively studied and identified regulators of drought response in our previous works [42,43]. This is why in our current study we shift our attention to saline stress in rice. We are specifically interested in observing if the *LYSA* regulators change under saline stress.

2. Materials and Methods

GRNs describe the complex interactions taking place between regulators and their target genes. Typically regulators consist of transcription factors (TFs), genes, RNA binding proteins, and regulator RNAs that can control the gene expression of the target genes [50–52]. GRNs govern the decision-making process in response to endogenous and external stimuli; thus, understanding their behavior at the genomic level can give us critical insights into achieving desirable phenotypical traits like increased lysine content [53,54].

GRNs have been modeled extensively in the past for a wide range of applications such as discovering novel biological relationships, studying complex diseases, drug design, and developing pathogen-resistant crops [55–59]. Common modeling techniques include differential equations, linear models, Boolean networks, probabilistic Boolean networks, Bayesian networks, and small molecule level models [60–64]. Each technique has its set of advantages and limitations. Therefore, we must consider the nature of the interactions in the GRN and the overall domain of the study while selecting a modeling method. In this paper, we are interested in modeling the lysine biosynthesis pathway in rice under normal (unstressed) and saline stress conditions. The interactions taking place in the pathway are sparse, multivariate, and stochastic in nature. Furthermore with the advent of high throughput technologies, publicly available genomic data is have become easily accessible [65]. Due to these factors, we will model the lysine biosynthesis pathway using Bayesian networks (BNs). BNs provide a stochastic framework and allow integration of pathway knowledge and data.

2.1. Bayesian Network Modeling

BNs are a class of Probabilistic Graphical Models (PGM) that integrate probability and graph theory to represent stochastic and causal relationships among variables in a system [66,67]. BNs consist of two main components (i) a directed acyclic graph (DAG) and (ii) local probability distributions (LPD) or the network parameters [68]. The DAG is a map that describes the causal relationships among the system variables, also known as nodes. DAGs specify the dependencies among the nodes and explain the flow of cause and effect in the overall network. The DAG can be derived from the literature or estimated from data using structure learning algorithms [69]. Associated with each node in the DAG is a local probability distribution (LPD) which describes the stochastic nature of interaction among the connected nodes [67]. The LPDs and the DAGs together describe the factorization of the joint probability distribution of all the nodes in terms of their LPDs. In order to formalize this notion consider a BN with N nodes such that it has a DAG structure $\mathcal{G}(X, E)$, where X_i represents the i^{th} node in the set of nodes X and E represents the set of casual edges between the nodes. Now suppose the LPD for each node X_i is given by $P(X_i | P_a(X_i))$, where $P_a(X_i)$ is the set of parent nodes of X_i . Then by the local Markov independence assumption, each node given its parent nodes, is independent of its nondescendant nodes. We can then factorize the joint probability of all the nodes in X as:

$$P(X = \{X_1, X_2, \dots, X_i, \dots, X_N\}) = \prod_{i=1}^N P(X_i | P_a(X_i)) \quad (1)$$

To model the lysine biosynthesis pathway using BN, we construct a DAG from the Kegg pathway we discussed in Figure 1. Learning the DAG from data is an NP-Hard problem and often requires selecting a graph structure from a candidate of possible DAGs [70,71]. This is a computationally expensive task, and the size of publicly available genomic datasets is not sufficiently large to produce a reliable DAG. Therefore we use pathway information (see Figure 1) to construct the DAG for the lysine biosynthesis pathway in Figure 2. Every node (represented by circles) in the DAG represents a gene present in the lysine biosynthesis pathway. These genes are referenced by their aliases; for instance, gene N represents *LYSA*. The nodes are connected by arrows that represent actual biological relationships as described in the pathway. We assume that genes in the network can be active, dormant, or inhibited. Thus we model each node as a categorical random variable with three states 1(active), 0 (dormant), and -1 (inhibited). Associated with each node is a rectangular box that describes the LPD (network parameter). For Node A , θ_A is vector representing the marginal probability of gene A being active, dormant, or inhibited. Similarly, $\theta_{M|L,K}$ is a vector representing the conditional probability of gene M being active, dormant, or inhibited given the states of its parents, gene L and gene K . This completes our discussion of the DAG for the lysine biosynthesis pathway. In the next section, we will

discuss how to estimate the LPDs. Once all the LPDs have been calculated the Bayesian network model is complete and can be used to perform gene intervention simulation under normal and saline stress conditions. These simulations will help us gain insight into the effect of intervening at the various genes. Genes that upregulate *LYSA* (gene N) will be considered ideal targets for genetic intervention. Interventions in the GRN can be carried out using gene editing methods such as CRISPR-Cas9 [54].

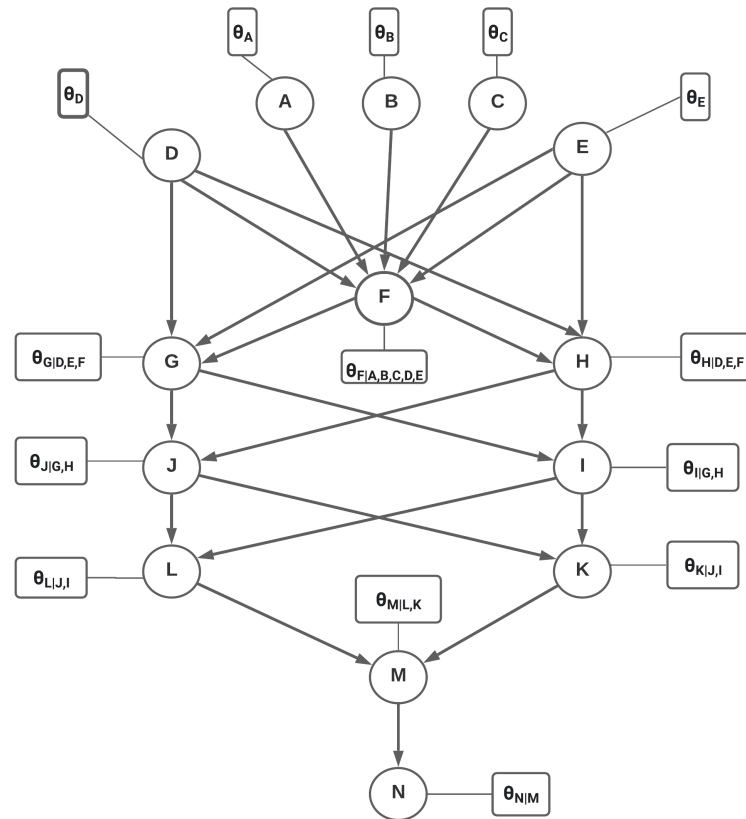


Figure 2. Directed Acyclic Graph (DAG) of the lysine biosynthesis pathway. Each node (circle) represents a gene in the pathway. The rectangular boxes represent the local probability distributions of the respective nodes. Each node is modeled as a categorical random variable with the following states: active(1), dormant(0), and inhibited (-1).

2.2. Parameter Estimation

Several methods can be employed to estimate the LPDs (network parameters) in a BN. Frequentist approaches such as Maximum Likelihood Estimation (MLE) are common when estimating the LPDs in a BN [72]. However, we will use a Bayesian approach to estimate the LPDs for the DAG constructed in the previous section. This is because the sizes of publicly available datasets are not sufficiently large to be reliably used by data-driven frequentist approaches. Unlike frequentist approaches, Bayesian estimation produces a posterior probability distribution for the LPDs based on data and prior knowledge [73]. The point estimate for the LPDs can be obtained by approximating the posterior distributions by their expected value or mode [74]. The Bayesian estimation process is based on Bayes rule where the posterior distribution of a random variable X , for a dataset \mathcal{D} , is given by:

$$P(X|\mathcal{D}) = \frac{P(\mathcal{D}|X)P(X)}{P(\mathcal{D})} \quad (2)$$

where $P(X)$ is the prior distribution of X

We will now use this approach to derive the general expression for estimating the LPDs for a BN where the nodes are modeled as categorical random variables. We can then extend our findings to the DAG in Figure 2.

Consider a BN with a DAG denoted by \mathcal{G} containing N (N is a Natural number) nodes. Each node X_i in \mathcal{G} is modeled as a categorical random variable with the following states: active (1), dormant (0), and inhibited (0). Thus for any node X_i in \mathcal{G} , $X_i \in \mathbf{S}=\{1,0,-1\}$, so if $X_i=0$, it implies that the node X_i is dormant. Let the probability with which X_i assumes any of the states in set \mathbf{S} be given by the probability vector θ_{X_i} . Then θ_{X_i} is of the form $[\theta_{X_i=1}, \theta_{X_i=0}, \theta_{X_i=-1}]^T$, where $\theta_{X_i=s}$ represents the probability of $X_i=s$ for $s \in \mathbf{S}$ and $\sum_s \theta_{X_i=s} = 1$. Now, suppose we have a dataset \mathcal{D} which contains n (n is natural number) independent and identically distributed (i.i.d) observations for each of the N nodes in \mathcal{G} . For a node X_i in \mathcal{G} , let $M_{X_i}[\mathbf{S}=s]$ represent the frequency of $X_i=s$ in \mathcal{D} ($\sum_s M_{X_i}[s]=n$). Then the likelihood under the dataset \mathcal{D} can be modeled as:

$$P(X_i | P_a(X_i), \theta_{X_i}) \sim \text{Multinomial}(\theta_{X_i}, n) \quad (3)$$

$$\text{Multinomial}(\theta_{X_i}, n) = n! \prod_{s \in \mathbf{S}} \frac{\theta_{X_i}^{M_{X_i}[s]}}{M_{X_i}[s]!} \quad (4)$$

The Bayesian estimation process requires selecting a prior distribution. Prior distributions can be selected based on domain knowledge; however in its absence, there are no fixed methods to choose a prior. The subjective selection of the prior distribution is often cited as a drawback of the Bayesian estimation process, as different priors lead to different results for the posterior distribution [75]. We set the prior distribution on θ_{X_i} for each node $X_i \in \mathcal{G}$ to follow a Dirichlet distribution. A Dirichlet prior under a multinomial likelihood causes the posterior distribution also to follow a Dirichlet distribution. This is because the multinomial and Dirichlet distributions belong to conjugate families of distributions [76,77]. Therefore we have the following formulation for the posterior distribution on θ_{X_i} :

$$\theta_{X_i} \sim \text{Dirichlet}(\boldsymbol{\alpha}) \quad (5)$$

$$\boldsymbol{\alpha} = [\alpha_{s=1}, \alpha_{s=0}, \alpha_{s=-1}]$$

$$\text{Dirichlet}(\theta_{X_i}; \boldsymbol{\alpha}) = \frac{1}{\beta(\boldsymbol{\alpha})} \prod_{s \in \mathbf{S}} [\theta_{X_i=s}]^{\alpha_s-1} \quad (6)$$

where $\beta(\boldsymbol{\alpha})$ is the Multivariate Beta function

$$P(\theta_{X_i} | X_i) = \text{Dirichlet}(\boldsymbol{\alpha}') \quad (7)$$

and

$$\boldsymbol{\alpha}' = [\alpha_{s=1} + M_{X_i}[s = 1], \alpha_{s=0} + M_{X_i}[s = 0], \alpha_{s=-1} + M_{X_i}[s = -1]]$$

$$\boldsymbol{\alpha}' = [\alpha'_{s=1}, \alpha'_{s=0}, \alpha'_{s=-1}]$$

In our study we specifically set the prior distribution on each node X_i to be Dirichlet($\alpha_{s=1}=1, \alpha_{s=0}=1, \alpha_{s=-1}=1$), which corresponds to uniform distribution over the open standard 2-simplex and is a non informative prior [78,79]. This is an appropriate choice for the prior distribution in our study as we do not have prior knowledge regarding the distribution of each node in the BN. Furthermore, this assumption on the prior distribution of the nodes allows us to obtain a closed form solution for the posterior distribution. Selecting a different prior will often lead to non-closed form solution for the posterior distribution and calculating the probability of data ($P(\mathcal{D})$) can be computationally expensive [80]. The formulation in Equation (7) represents the posterior distribution of the node parameter θ_{X_i} .

We approximate θ_{X_i} by its expected value in order to obtain a point estimate for the LPDs in the BN. The expectation of a Dirichlet distribution is given by [81]:

$$\theta_{X_i} = \begin{bmatrix} \theta_{X_i=1} \\ \theta_{X_i=0} \\ \theta_{X_i=-1} \end{bmatrix} \approx E[\theta_{X_i}|X_i] = \begin{bmatrix} \frac{\alpha'_{s=1}}{\sum_S \alpha'_s} \\ \frac{\alpha'_{s=0}}{\sum_S \alpha'_s} \\ \frac{\alpha'_{s=-1}}{\sum_S \alpha'_s} \end{bmatrix} \quad (8)$$

Similarly if we have a node X_i with a parent node $Y_i=s$ ($s \in \mathbf{S}$) under the same Dirichlet and Multinomial framework, then the LPD associated with $\theta_{X_i|Y_i}$ can be formulated as follows:

$$\theta_{X_i|Y_i=s} = \begin{bmatrix} \theta_{X_i=1|Y_i=s} \\ \theta_{X_i=0|Y_i=s} \\ \theta_{X_i=-1|Y_i=s} \end{bmatrix} \approx E[\theta_{X_i|Y_i=s} | (X_i | Y_i = s)] = \begin{bmatrix} \frac{\alpha_{s=1} + M_{X_i|Y_i}[X_i=1, Y_i=s]}{\sum_S \alpha_s + M_{X_i|Y_i}[X_i=1, Y_i=s]} \\ \frac{\alpha_{s=0} + M_{X_i|Y_i}[X_i=0, Y_i=s]}{\sum_S \alpha_s + M_{X_i|Y_i}[X_i=0, Y_i=s]} \\ \frac{\alpha_{s=-1} + M_{X_i|Y_i}[X_i=-1, Y_i=s]}{\sum_S \alpha_s + M_{X_i|Y_i}[X_i=-1, Y_i=s]} \end{bmatrix} \quad (9)$$

In equation (9), $M_{X_i|Y_i}[X_i=1, Y_i=s]$ represents the frequencies when $X_i=1$ and $Y_i=s$ simultaneously in the dataset \mathcal{D} . Similarly, $M_{X_i|Y_i}[X_i=0, Y_i=s]$ is the frequency of datapoints in \mathcal{D} when $X_i=0$ and $Y_i=s$ simultaneously, and so on for $X_i=-1$. Once the node parameters are estimated, gene intervention simulations can be carried out using inference in the BN. Inference computes the effect of intervening at each node on the reporter gene *LYSA* (gene N).

2.3. Gene Intervention Simulations

BNs represent the cause and effect relationship among the nodes of the system being modeled. Inference quantifies the cause and effect relationship by allowing us to compute conditional probability queries. Then for a node of interest X , also known as the query node and an intervention (or evidence) node E in the BN, we can compute the conditional probability $P(X|E)$ using inference algorithms. This implies, if we instantiate (fix) node E , we can calculate its effect on node X . Inference algorithms use the network parameters and structural dependencies to compute the required conditional probabilities. To further elucidate this notion, consider the BN shown in Figure 3. Let each node of the BN be a binary random variable with states 0 and 1. Suppose we have estimated the LPDs $P(A)$, $P(B|A)$, $P(C|A)$, $P(D|B,C)$, then we can use inference in this BN to answer conditional probability queries such as $P(D=1|A=1)$.

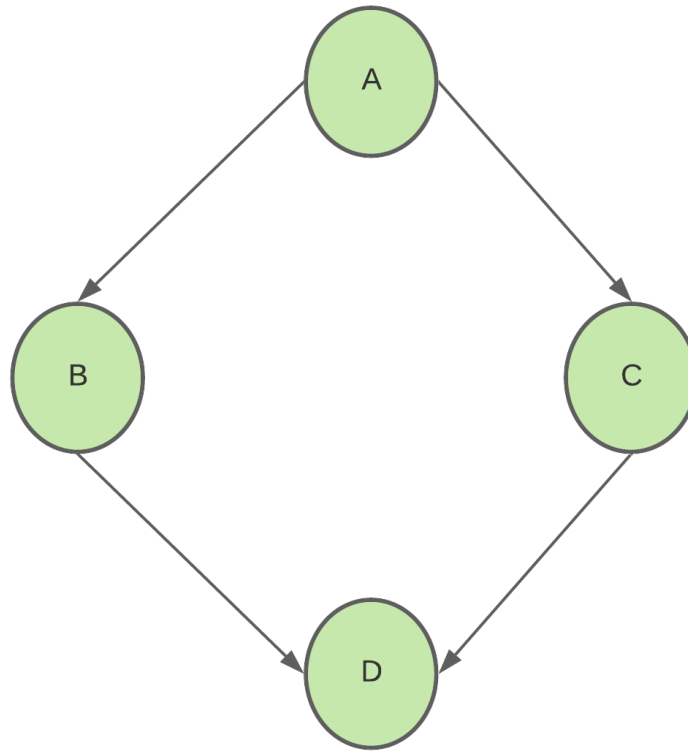


Figure 3. Example BN with binary nodes.

We compute $P(D=1 | A=1)$ as follows:

$$\begin{aligned}
 P(D = 1 | A = 1) &= \frac{P(D = 1, A = 1)}{P(A = 1)} \\
 &= \frac{\sum_B \sum_C P(A = 1, B, C, D = 1)}{P(A = 1)}
 \end{aligned}$$

Using the properties of the BN, all nodes are independent of any non descendant nodes

$$\begin{aligned}
 P(D = 1 | A = 1) &= \frac{\sum_B \sum_C P(A = 1)P(B|A)P(C|A = 1)P(D = 1|B, C)}{P(A = 1)} \\
 &= \sum_B \sum_C P(B|A = 1)P(C|A = 1)P(D = 1|B, C)
 \end{aligned}$$

We can use the LPDs to calculate the exact probability $P(D=1 | A=1)$.

Inference techniques such as the one applied in the BN in Figure 3 are classified as "exact" because they compute the true values for the conditional probability query. However, exact inference in BNs has been shown to be NP-hard [82,83]. While there exist efficient algorithms for exact inference, they are often limited to simpler DAG structures[82]. For example, Pearl's message-passing algorithm works efficiently for singly connected DAG structures [84]. Therefore for larger DAGs, exact inference is not ideal as the computational cost of calculating the conditional probabilities can be expensive. In such cases, we employ approximate inference algorithms, which produce estimates of the exact conditional probabilities [85]. Approximate inference can include wide-ranging techniques such as model simplification methods, loopy belief propagation methods, search based methods, utility based methods, and stochastic simulation methods [86]. In this paper, we implement a stochastic simulation-based inference technique called Likelihood Weighting (LW) to estimate the conditional probability queries in the BN model for the lysine Biosynthesis

pathway. Stochastic simulation techniques estimate the conditional probabilities by drawing samples from the LPDs. These estimates typically converge to the true conditional probabilities as the number of samples drawn increases. LW can efficiently handle inference of large multiply connected BNs and is based on forward sampling [86,87]. Since our BN model is multiply connected and we are only interested in estimating $P(N=1 \mid E \in \{A,B,C,\dots,M\})$, i.e., the probability of upregulating *LYSA* (gene *N*), while conditioning on other genes (evidence or intervention nodes), LW turns out to be a suitable method for performing inference.

The LW algorithm estimates the conditional probability, $P(X=x \mid E=e)$ for a query node *X* and an evidence node *E* by generating samples from a BN model. We fix the sample size (*m*) and a topological ordering at the start of the algorithm. The algorithm iterates through a sample generation process *m* times, and then computes the conditional probability from the generated samples. During the sample generation process, the algorithm generates values for the nonevidence nodes only; it sets the value of the evidence node to its observed (*e*, in this case) value. The node values for each sample are generated in the established topological ordering. Each sample is assigned a weight of 1 at the start of the sample generation process. The weight is updated only when an evidence node is encountered while traversing the topological ordering. When this happens, the sample's weight is updated by multiplying the current weight with the likelihood of the evidence node conditioned on the state of its parent nodes. The likelihood is given by the probability $P(E=e \mid P_n(E))$. The process is repeated until *m* samples are generated. Following this step, conditional probability is estimated by dividing the sum of the weights of the samples where $X=x$ by the sum of all the sample's weights. The pseudo code for the LW algorithm by Stuart Russell and Peter Norvig is presented in Algorithm 1 [88].

Algorithm 1: Likelihood-Weighting Algorithm

```

Function LIKELIHOOD-WEIGHTING(X, e, bn, N):
  outputs an estimate of  $P(X \mid e)$ 
  inputs: X, the query variable
           e, observed values for variables E
           bn, a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
           N, the total number of samples to be generated
  local variables: W, a vector of weighted counts for each value of X, initially zero
  for j=1 to N do
    x, w ← WEIGHTED-SAMPLE(bn, e)
    W[x] ← W[x] + w where x is the value of X in x
  end
  return NORMALIZE(W)

Function WEIGHTED-SAMPLE(bn, e):
  outputs an event and a weight
  w ← 1; x ← an event with n elements initialized from e
  for each variable Xi in X1, ..., Xn do
    if Xi is an evidence variable with value xi in e then
      w ← w ×  $P(X_i = x_i \mid \text{parents}(i))$ 
    else
      x[i] ← a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
    end
  end
  return x, w

```

2.4. Dataset

To estimate the LPDs in the BN model, we use the dataset GSE98455, which is publicly available from the NCBI GEO database [89–91]. This dataset contains RNA-Seq counts for

rice seedlings under saline stress and normal (unstressed or control) conditions. Since our BN model contains nodes modeled as categorical variables, the RNA-Seq data had to be preprocessed. The first preprocessing step was to normalize the count data. This was done using the ratio of medians methods described in the DESeq2 data processing protocols by *Love et al.*[92]. DESeq2 is one of the most commonly used RNA-Seq data processing protocols and is easily accessible on the R programming language as a package (DESeq2)[93–96]. After this step, we extracted the normalized data for genes in our BN model from the entire dataset. We further segregated the normalized dataset based on saline stress and normal conditions. The two resulting datasets contained normalized data only for Genes A-N under saline stress and normal conditions. Each of the normalized datasets was then transformed to categorical values using K-means clustering. The clustering process categorized the data in both the datasets into the following values 1 (active), 0 (dormant), and -1 (inhibited). Both the unstressed (normal) and saline stressed datasets had 184 datapoints per gene, and the LPDs were computed for each of these datasets. A visual representation of the data processing pipeline has been presented in Figure 4. Figures 5 and 6 show the discretized categorical data for each node in the BN under normal and saline stress conditions, respectively.

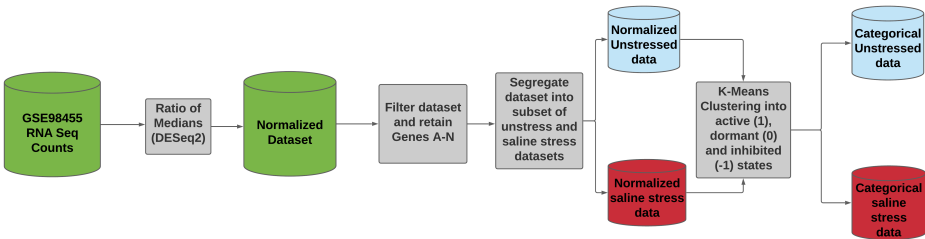


Figure 4. Data processing pipeline for RNA-Seq dataset GSE98455.

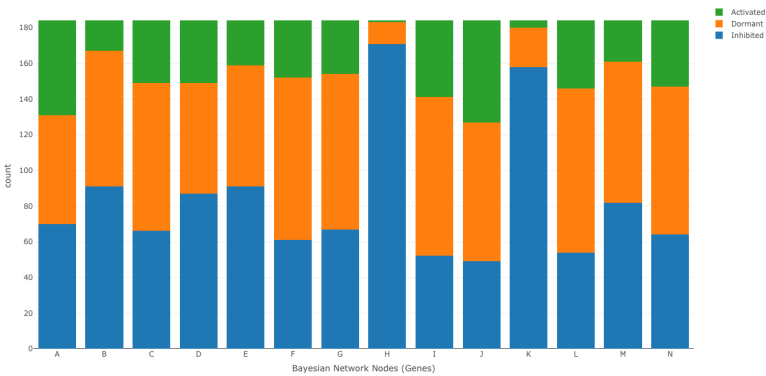


Figure 5. Discretized RNA-Seq data under normal conditions.

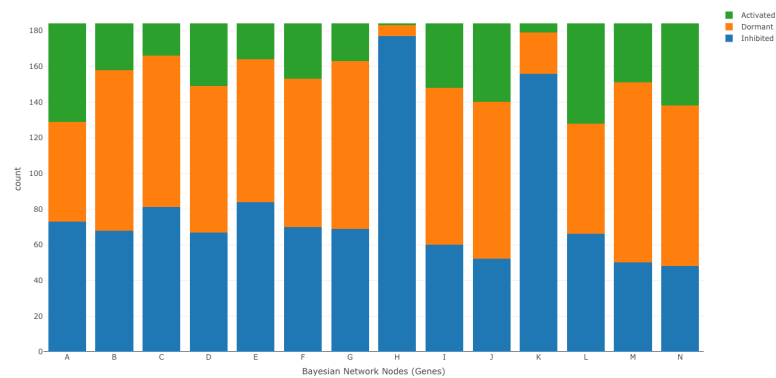


Figure 6. Discretized RNA-Seq data under saline stress conditions.

3. Results

The LPDs estimated from the RNA-Seq dataset were used to simulate gene intervention in the BN. When intervening at a gene, the node representing that gene in the BN was instantiated to a status of active (1), dormant (0), or inhibited (-1). We applied the LW algorithm with a large sample size of 600,000 to compute the probability $P(N=1 \mid \text{Gene Intervention})$ and ensure convergence of the probabilities being estimated. Gene N (*LYSA*) is set as the query node because it is the reporter gene for lysine production, thus upregulating gene N (*LYSA*) may lead to increased lysine production. We perform intervention at genes A-M one at a time and then in combinations of two (pairs) at a time. These gene intervention strategies were applied under both normal and saline stress conditions. In order to measure the causal effect of intervention, we subtract the marginal probability $P(N=1)$ from $P(N=1 \mid \text{Gene Intervention})$, for all the possible gene intervention strategies. This difference is defined as the score metric and is used to compare the effectiveness of each gene intervention strategy. The data processing and probability computation pipeline was written in the R programming language, and the Bnlearn package was used to perform LW [97–99]. So,

$$\text{score} = P(N = 1 \mid \text{Gene Intervention}) - P(N = 1). \quad (10)$$

Since there are many possible combinations under single and pairwise gene interventions, we only include the top five intervention strategies with the highest scores in Figures 7-8. In Figures 7.a and 7.b, we present the scores for single node intervention under normal and saline stress conditions.

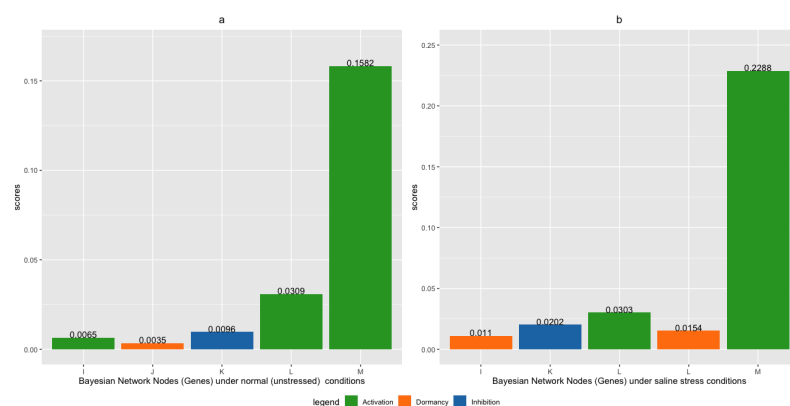


Figure 7. Single node intervention under (a) normal and (b) saline stress conditions.

It is clear from Figures 7.7 (a) and 7.7 (b) that activating gene M (*DAPF*) has the maximum score. This implies that under both normal and saline stress conditions, genetically activating gene M (*DAPF*) has the best chance for upregulating the reporter gene N

(*LYSA*). We also notice that gene L (*AGD2*) is also fairly active in its role in upregulating gene N (*LYSA*). Activating gene L (*AGD2*) achieves the second-largest score under normal conditions. Under saline stress conditions, activating gene L (*AGD2*) or keeping it dormant also ranks among the top five gene intervention strategies. Inhibiting gene K (*ALD1*) achieves the third and second-highest scores under normal and saline stress conditions. Additionally, we also observe that midstream genes such as gene I (*DAPB1*) and gene J (*DAPB2*) also play an active role in upregulating gene N (*LYSA*). However, activating gene M (*DAPF*) has a significantly higher score under both conditions, thus gene M (*DAPF*) serves as an ideal candidate for gene intervention.

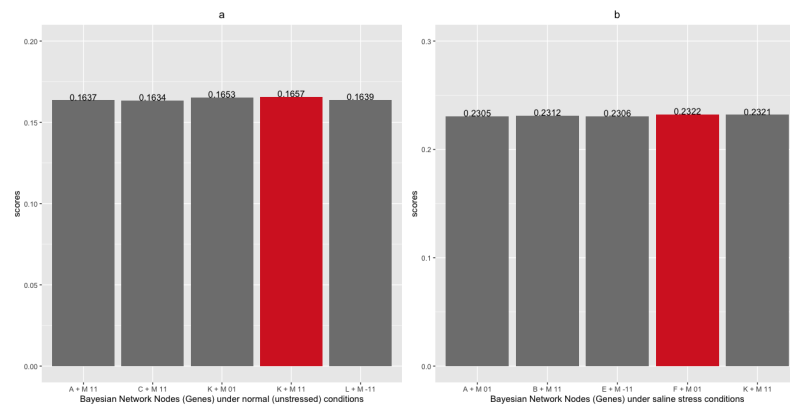


Figure 8. Pairwise node intervention under (a) normal and (b) saline stress conditions.

Figures 8.8 (a) and (b) represent the five highest-scoring pairwise intervention strategies for normal and saline stress conditions. Under each condition, the score amongst different strategies are almost similar with marginal differences. We highlight the highest-scoring strategy in red. Under normal conditions, activating both gene K (*ALD1*) and gene M (*DAPF*) maximized the scores. While under saline stress keeping gene F (*LOC_Os03g55280*) dormant and activating gene M (*DAPF*) achieved the highest score. This implies that under each of the conditions, the respective pairwise intervention strategy with highest scores maximize the likelihood of upregulating gene N (*LYSA*). In Figures 8.8 (a) and (b), it can also be seen that upstream genes such as genes A,B,C and E are also involved in the upregulation of gene N (*LYSA*) and produce comparable scores to those produced by the regulation of downstream genes like gene K (*ALD1*) and gene L (*AGD2*). Across both the conditions, we also observe that gene M (*DAPF*) is always upregulated, which serves to be a further indicator of the high regulatory effect of gene M (*DAPF*) on gene N (*LYSA*).

4. Discussion

In this paper, we studied the lysine biosynthesis pathway in rice to identify the genetic regulators of lysine content. Rice is a staple food source for 50% of the global population; with lysine being the first limiting essential amino acid in rice, it is vital to identify gene regulators that can boost lysine content. We modeled the lysine biosynthesis pathway in rice using BNs under normal and saline stress conditions to identify these regulators. We used BNs because they allow us to integrate domain knowledge in the form of pathway information with experimental data. We used real-world RNA-Seq data to estimate the LPDs in the BN and run the gene intervention simulations. We intervened at the genes one at a time and then in pairwise combinations using the LW inference algorithm. We calculated a score metric to measure the efficacy of the gene intervention strategies.

Our analysis revealed that upregulating *DAPF* (gene M) maximized the probability of the lysine reporter gene *LYSA* (gene N) being upregulated under both normal and saline stress conditions. When *DAPF* (gene M) was upregulated, it not only achieved the highest score under single gene intervention but was also present in all the five highest-scoring gene intervention strategies under pairwise intervention. This implies that *DAPF* (gene M) is a positive regulator of *LYSA* (gene N) and serves as an ideal candidate for genetic

intervention. Gene editing can be used to target and upregulate *DAPF* (gene M) in rice. Field experiments involving *DAPF* overexpressing rice can confirm if this intervention strategy upregulates *LYSA* and increases the overall lysine content. We further observed under single gene intervention that midstream genes such as *DAPB1* (gene I) and *DAPB2* (gene J) also played significant role in upregulating *LYSA* (gene N). On the other hand, under pairwise intervention, we found upstream genes such as genes A, B, C, and E were also involved in upregulating *LYSA* (gene N).

The future steps in our study of lysine will include confirming our finding in this paper by performing validation experiments in the field. We would also like to improve our choice of the prior distribution on each node. In our current analysis, we used a noninformative prior as we did not have any knowledge regarding the prior distribution of the nodes in the BN. Using informative priors may increase the computational costs but has the potential to improve our predictions of lysine regulators. Furthermore, we are also interested in studying how other essential amino acids such as Threonine, Methionine, and Isoleucine in the larger aspartate pathway regulate lysine content. Threonine is known to downregulate the enzyme AK in the lysine biosynthesis pathway; thus, studying the multilevel regulation among the different amino acids in the aspartate pathway will help understand lysine production.

5. Conclusions

We modeled the lysine biosynthesis pathway in rice under normal and saline stress conditions to identify the regulators of lysine. Among the essential amino acids, lysine is present in the least quantity in rice; thus increasing its content in rice will improve its nutritional value. Our analysis revealed that under both the normal and saline stress conditions upregulating *DAPF* is the best genetic intervention strategy for upregulating the lysine reporter gene *LYSA*. Applying gene intervention techniques such as CRISPR-Cas9 to upregulate *DAPF* has the potential to increase the lysine content in rice.

Author Contributions: Conceptualization, E.S. and A.D.; methodology, A.L. and K.R.; software, A.L.; validation E.S, K.R.; formal analysis, A.L. and A.D.; investigation, A.L. and A.D.; resources, A.D. and E.S.; data curation, A.L.; writing—original draft preparation, A.L.; writing—review and editing, A.L., A.D., K.R. and E.S.; visualization, A.L.; supervision, E.S. and A.D.; project administration, E.S. and A.D.; funding acquisition, E.S. and A.D.. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported in part by the TEES-AgriLife Center for Bioinformatics and Genomic Systems Engineering (CBGSE) startup funds, the Texas A&M X-Grant Program, and in part by the National Science Foundation under Grant ECCS-1609236(to A. D.). The funding bodies did not play any role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Institutional Review Board Statement: Not applicable

Informed Consent Statement: Not applicable

Data Availability Statement: The datasets used in this study are publicly available at the NCBI with the accession numbers of GSE98455. The subset of data extracted from the dataset to support the conclusion of this article are included within the article. The R code containing the simulations are provided as additional files.

Acknowledgments: We would like to recognize Dr. Chi Zhang, Associate Professor, School of Biological Sciences, University of Nebraska, Lincoln. He generously provided us with the gene annotation file required for analyzing the dataset GSE98455

Conflicts of Interest: The authors declare no conflict of interest.

Sample Availability: Not applicable

Abbreviations

The following abbreviations are used in this manuscript:

LKR	Lysine Ketoglutarate Reductase
SDH	Saccharopine Dehydrogenase
DHPS	Dihydrodipicolinate Synthase
AK	Aspartate Kinase
GRN	Gene Regulatory Network
GMO	genetically modified organisms
MSU	Michigan State University
TF	Transcription Factor
BN	Bayesian Network
PGM	Probabilistic Graphical Model
LPD	Local Probability Distribution
i.i.d	Independent and Identically Distributed
LW	Likelihood Weighting

References

1. Alberts, B.; Bray, D.; Hopkin, K.; Johnson, A.; Lewis, J.; Raff, M.; Roberts, K.; Walter, P. *Essential Cell Biology*, 3 ed.; Garland Science, 2010; pp. 119–122.

2. Alberts, B.; Johnson, A.; Lewis, J.; Raff, M.; Roberts, K.; Walter, P. The shape and structure of protein. In *Molecular Biology of the Cell*, 4 ed.; Garland Science, 2002.

3. Lopez, M.; Mohiuddin, S. Biochemistry, Essential Amino Acids. *StatPearls Publishing* **2020**.

4. D’Mello, J.P.F. Amino acids as multifunctional molecules. In *Amino Acids in Animal Nutrition*, 2 ed.; CABI Publishing, 2003; pp. 2–2.

5. Hoffman, J.; Falvo, M. Protein - Which is Best? *Journal of Sports Science and Medicine* **2004**.

6. Tien Lea, D.; Duc Chua, H.; Quynh Lea, N. Improving Nutritional Quality of Plant Proteins Through Genetic Engineering. *Current Genomics* **2016**, 17. doi:10.2174/1389202917666160202215934.

7. Stencel, C.; Dobbins, C. Report Offers New Eating and Physical Activity Targets To Reduce Chronic Disease Risk, 2002.

8. Zha, Y.; Qian, Q. Protein Nutrition and Malnutrition in CKD and ESRD. *Nutrients* **2017**, 9. doi:10.3390/nu9030208.

9. National Research Council. *Recommended Dietary Allowances*, 10 ed.; National Academies Press: Washington, D.C., 1989. doi:10.17226/1349.

10. Titchenal, A.; Hara, S.; Arceo Caacbay, N.; Meinke-Lau, W.; Yang, Y.Y.; Ksinoa Fialkowski Revilla, M.; Draper, J.; Langfelder, G.; Gibby, C.; Nicole Chun, C.; Calabrese, A. *Human Nutrition*, 2020 ed.; University of Hawaii at Manoa Food Science and Human Nutrition Program, 2020; pp. 395–402.

11. Henchion, M.; Hayes, M.; Mullen, A.; Fenelon, M.; Tiwari, B. Future Protein Supply and Demand: Strategies and Factors Influencing a Sustainable Equilibrium. *Foods* **2017**, 6. doi:10.3390/foods6070053.

12. Vasileška, A.; Rechkoska, G. Global and Regional Food Consumption Patterns and Trends. *Procedia - Social and Behavioral Sciences* **2012**, 44. doi:10.1016/j.sbspro.2012.05.040.

13. Berrazaga, I.; Micard, V.; Gueugneau, M.; Walrand, S. The Role of the Anabolic Properties of Plant- versus Animal-Based Protein Sources in Supporting Muscle Mass Maintenance: A Critical Review. *Nutrients* **2019**, 11. doi:10.3390/nu11081825.

14. de Gavelle, E.; Huneau, J.F.; Bianchi, C.; Verger, E.; Mariotti, F. Protein Adequacy Is Primarily a Matter of Protein Quantity, Not Quality: Modeling an Increase in Plant:Animal Protein Ratio in French Adults. *Nutrients* **2017**, 9. doi:10.3390/nu9121333.

15. Abete, I.; Romaguera, D.; Vieira, A.R.; Lopez de Munain, A.; Norat, T. Association between total, processed, red and white meat consumption and all-cause, CVD and IHD mortality: a meta-analysis of cohort studies. *British Journal of Nutrition* **2014**, 112. doi:10.1017/S000711451400124X.

16. Demeyer, D.; Mertens, B.; De Smet, S.; Ulens, M. Mechanisms Linking Colorectal Cancer to the Consumption of (Processed) Red Meat: A Review. *Critical Reviews in Food Science and Nutrition* **2016**, 56. doi:10.1080/10408398.2013.873886.

17. Malik, V.S.; Li, Y.; Tobias, D.K.; Pan, A.; Hu, F.B. Dietary Protein Intake and Risk of Type 2 Diabetes in US Men and Women. *American Journal of Epidemiology* **2016**, 183. doi:10.1093/aje/kwv268.

18. The Food and Agriculture Organization of the United Nations. Livestock solutions for climate change. Technical report, United Nations, 2017.

19. United Nations. World population projected to reach 9.8 billion in 2050, and 11.2 billion in 2100, 2017.

20. Day, L. Proteins from land plants – Potential resources for human nutrition and food security. *Trends in Food Science & Technology* **2013**, 32. doi:10.1016/j.tifs.2013.05.005.

21. Rosegrant, M.W.; Leach, N.; Gerpacio, R.V. Alternative futures for world cereal and meat consumption. *Proceedings of the Nutrition Society* **1999**, 58. doi:10.1017/S0029665199000312.

22. Millward, D.J.; Jackson, A.A. Protein/energy ratios of current diets in developed and developing countries compared with a safe protein/energy ratio: implications for recommended protein and amino acid intakes. *Public Health Nutrition* **2004**, 7. doi:10.1079/PHN2003545.

23. Kusano, M.; Yang, Z.; Okazaki, Y.; Nakabayashi, R.; Fukushima, A.; Saito, K. Using Metabolomic Approaches to Explore Chemical Diversity in Rice. *Molecular Plant* **2015**, *8*. doi:10.1016/j.molp.2014.11.010.
24. Galili, G.; Amir, R. Fortifying plants with the essential amino acids lysine and methionine to improve nutritional quality. *Plant Biotechnology Journal* **2013**, *11*. doi:10.1111/pbi.12025.
25. Wang, W.; Galili, G. Transgenic high-lysine rice – a realistic solution to malnutrition? *Journal of Experimental Botany* **2016**, *67*. doi:10.1093/jxb/erw254.
26. Galili, G.; Karchi, H.; Shaul, O.; Perl, A.; Cahana, A.; Tzchori, I.B.T.; Zhu, X.Z.; Galili, S. Production of transgenic plants containing elevated levels of lysine and threonine. *Biochemical Society Transactions* **1994**, *22*. doi:10.1042/bst0220921.
27. Grigg, D. The pattern of world protein consumption. *Geoforum* **1995**, *26*. doi:10.1016/0016-7185(94)00020-8.
28. Juliano, B.O.; The Food and Agriculture Organization of the United Nations. World rice production compared to other cereals. In *Rice in human nutrition*; International Rice Research Institute of the United Nations: Rome, 1993.
29. Muthayya, S.; Sugimoto, J.D.; Montgomery, S.; Maberly, G.F. An overview of global rice production, supply, trade, and consumption. *Annals of the New York Academy of Sciences* **2014**, *1324*. doi:10.1111/nyas.12540.
30. Kawakatsu, T.; Takaiwa, F. Differences in Transcriptional Regulatory Mechanisms Functioning for Free Lysine Content and Seed Storage Protein Accumulation in Rice Grain. *Plant and Cell Physiology* **2010**, *51*. doi:10.1093/pcp/pcq164.
31. Arruda, P.; Kemper, E.L.; Papes, F.; Leite, A. Regulation of lysine catabolism in higher plants. *Trends in Plant Science* **2000**, *5*. doi:10.1016/S1360-1385(00)01688-5.
32. Long, X.; Liu, Q.; Chan, M.; Wang, Q.; Sun, S.S.M. Metabolic engineering and profiling of rice with increased lysine. *Plant Biotechnology Journal* **2013**, *11*. doi:10.1111/pbi.12037.
33. Yang, Q.q.; Zhang, C.q.; Chan, M.I.; Zhao, D.s.; Chen, J.z.; Wang, Q.; Li, Q.f.; Yu, H.x.; Gu, M.h.; Sun, S.S.m.; Liu, Q.q. Biofortification of rice with the essential amino acid lysine: molecular characterization, nutritional evaluation, and field performance. *Journal of Experimental Botany* **2016**, *67*. doi:10.1093/jxb/erw209.
34. Zhu, X.; Galili, G. Increased Lysine Synthesis Coupled with a Knockout of Its Catabolism Synergistically Boosts Lysine Content and Also Transregulates the Metabolism of Other Amino Acids in Arabidopsis Seeds. *The Plant Cell* **2003**, *15*. doi:10.1105/tpc.009647.
35. Angelovici, R.; Fait, A.; Fernie, A.R.; Galili, G. A seed high-lysine trait is negatively associated with the TCA cycle and slows down Arabidopsis seed germination. *New Phytologist* **2011**, *189*. doi:10.1111/j.1469-8137.2010.03478.x.
36. Tzchori, I.B.T.; Perl, A.; Galili, G. Lysine and threonine metabolism are subject to complex patterns of regulation in Arabidopsis. *Plant Molecular Biology* **1996**, *32*. doi:10.1007/BF00020213.
37. Rappe, M. CRISPR Plants: New Non-GMO Method to Edit Plants, 2020.
38. Shew, A.M.; Nalley, L.L.; Snell, H.A.; Nayga, R.M.; Dixon, B.L. CRISPR versus GMOs: Public acceptance and valuation. *Global Food Security* **2018**, *19*, 71–80. doi:10.1016/j.gfs.2018.10.005.
39. Rastogi, K.; Ibarra, O.; Molina, M.; Faion-Molina, M.; Thomson, M.; Septiningsih, E.M. Using CRISPR/Cas9 Genome Editing to Increase Lysine Levels in Rice. ASA-CSSA-SSA International Annual Meeting, San Antonio, TX; , 2019.
40. Kanehisa, M. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Research* **2000**, *28*. doi:10.1093/nar/28.1.27.
41. Kawahara, Y.; de la Bastide, M.; Hamilton, J.P.; Kanamori, H.; McCombie, W.R.; Ouyang, S.; Schwartz, D.C.; Tanaka, T.; Wu, J.; Zhou, S.; Childs, K.L.; Davidson, R.M.; Lin, H.; Quesada-Ocampo, L.; Vaillancourt, B.; Sakai, H.; Lee, S.S.; Kim, J.; Numa, H.; Itoh, T.; Buell, C.R.; Matsumoto, T. Improvement of the *Oryza sativa* Nipponbare reference genome using next generation sequence and optical map data. *Rice* **2013**, *6*. doi:10.1186/1939-8433-6-4.
42. Lahiri, A.; Venkatasubramani, P.S.; Datta, A. Bayesian modeling of plant drought resistance pathway. *BMC Plant Biology* **2019**, *19*. doi:10.1186/s12870-019-1684-3.
43. Lahiri, A.; Zhou, L.; He, P.; Datta, A. Detecting Drought Regulators using Stochastic Inference in Bayesian Networks. Manuscript submitted for publication. doi:10.21203/rs.3.rs-73056/v1.
44. Sheng, M.; Tang, M.; Chen, H.; Yang, B.; Zhang, F.; Huang, Y. Influence of arbuscular mycorrhizae on photosynthesis and water status of maize plants under salt stress. *Mycorrhiza* **2008**, *18*. doi:10.1007/s00572-008-0180-7.
45. Tisarum, R.; Theerawitaya, C.; Samphumphuang, T.; Polispitak, K.; Thongpoem, P.; Singh, H.P.; Cha-um, S. Alleviation of Salt Stress in Upland Rice (*Oryza sativa* L. ssp. indica cv. Leum Pua) Using Arbuscular Mycorrhizal Fungi Inoculation. *Frontiers in Plant Science* **2020**, *11*. doi:10.3389/fpls.2020.00348.
46. Reddy, I.N.B.L.; Kim, B.K.; Yoon, I.S.; Kim, K.H.; Kwon, T.R. Salt Tolerance in Rice: Focus on Mechanisms and Approaches. *Rice Science* **2017**, *24*. doi:10.1016/j.rsci.2016.09.004.
47. Kakar, N.; Jumaa, S.H.; Redoña, E.D.; Warburton, M.L.; Reddy, K.R. Evaluating rice for salinity using pot-culture provides a systematic tolerance assessment at the seedling stage. *Rice* **2019**, *12*. doi:10.1186/s12284-019-0317-7.
48. Deshmukh, V.; Mankar, S.P.; Muthukumar, C.; Divahar, P.; Bharathi, A.; Thomas, H.B.; Rajurkar, A.; Sellamuthu, R.; Poornima, R.; Senthivel, S.; et al.. Genome-Wide Consistent Molecular Markers Associated with Phenology, Plant Production and Root Traits in Diverse Rice (*Oryza sativa* L.) Accessions under Drought in Rainfed Target Populations of the Environment. *Current Science* **2018**, *114*, 329–340. doi:10.18520/cs/v114/i02/329-340.
49. Razzaque, S.; Elias, S.M.; Haque, T.; Biswas, S.; Jewel, G.M.N.A.; Rahman, S.; Weng, X.; Ismail, A.M.; Walia, H.; Juenger, T.E.; Seraj, Z.I. Gene Expression analysis associated with salt stress in a reciprocally crossed rice population. *Scientific Reports* **2019**, *9*. doi:10.1038/s41598-019-44757-4.

50. Jackson, C.A.; Castro, D.M.; Saldi, G.A.; Bonneau, R.; Gresham, D. Gene regulatory network reconstruction using single-cell RNA sequencing of barcoded genotypes in diverse environments. *eLife* **2020**, *9*. doi:10.7554/eLife.51254.
51. Davidson, E.H.; Erwin, D.H. Gene Regulatory Networks and the Evolution of Animal Body Plans. *Science* **2006**, *311*. doi:10.1126/science.1113832.
52. Kaern, M.; Blake, W.J.; Collins, J. The Engineering of Gene Regulatory Networks. *Annual Review of Biomedical Engineering* **2003**, *5*. doi:10.1146/annurev.bioeng.5.040202.121553.
53. Bonnaffoux, A.; Herbach, U.; Richard, A.; Guillemain, A.; Gonin-Giraud, S.; Gros, P.A.; Gandrillon, O. WASABI: a dynamic iterative framework for gene regulatory network inference. *BMC Bioinformatics* **2019**, *20*. doi:10.1186/s12859-019-2798-1.
54. Sun, Y.; Dinneny, J.R. Q&A: How do gene regulatory networks control environmental responses in plants? *BMC Biology* **2018**, *16*. doi:10.1186/s12915-018-0506-7.
55. Emmert-Streib, F.; Dehmer, M.; Haibe-Kains, B. Gene regulatory networks and their applications: understanding biological and medical problems in terms of networks. *Frontiers in Cell and Developmental Biology* **2014**, *2*. doi:10.3389/fcell.2014.00038.
56. Arshad, O.A.; Datta, A. Towards targeted combinatorial therapy design for the treatment of castration-resistant prostate cancer. *BMC Bioinformatics* **2017**, *18*. doi:10.1186/s12859-017-1522-2.
57. Vundavilli, H.; Datta, A.; Sima, C.; Hua, J.; Lopes, R.; Bittner, M. Targeting oncogenic mutations in colorectal cancer using cryptotanshinone. *PLOS ONE* **2021**, *16*. doi:10.1371/journal.pone.0247190.
58. Timmermann, T.; González, B.; Ruz, G.A. Reconstruction of a gene regulatory network of the induced systemic resistance defense response in Arabidopsis using boolean networks. *BMC Bioinformatics* **2020**, *21*. doi:10.1186/s12859-020-3472-3.
59. Venkat, P.S.; Narayanan, K.R.; Datta, A. A Bayesian Network-Based Approach to Selection of Intervention Points in the Mitogen-Activated Protein Kinase Plant Defense Response Pathway. *Journal of Computational Biology* **2017**, *24*. doi:10.1089/cmb.2016.0089.
60. Vijesh, N.; Chakrabarti, S.K.; Sreekumar, J. Modeling of gene regulatory networks: A review. *Journal of Biomedical Science and Engineering* **2013**, *06*. doi:10.4236/jbise.2013.62A027.
61. Vundavilli, H.; Datta, A.; Sima, C.; Hua, J.; Lopes, R.; Bittner, M. Bayesian Inference Identifies Combination Therapeutic Targets in Breast Cancer. *IEEE Transactions on Biomedical Engineering* **2019**, *66*. doi:10.1109/TBME.2019.2894980.
62. Karlebach, G.; Shamir, R. Modelling and analysis of gene regulatory networks. *Nature Reviews Molecular Cell Biology* **2008**, *9*. doi:10.1038/nrm2503.
63. Vundavilli, H.; Datta, A.; Sima, C.; Hua, J.; Lopes, R.; Bittner, M. Cryptotanshinone Induces Cell Death in Lung Cancer by Targeting Aberrant Feedback Loops. *IEEE Journal of Biomedical and Health Informatics* **2020**, *24*. doi:10.1109/JBHI.2019.2958042.
64. Kapoor, R.; Datta, A.; Sima, C.; Hua, J.; Lopes, R.; Bittner, M.L. A Gaussian Mixture-Model Exploiting Pathway Knowledge for Dissecting Cancer Heterogeneity. *IEEE/ACM Transactions on Computational Biology and Bioinformatics* **2019**, p. 1–1. doi:10.1109/TCBB.2018.2869813.
65. Sinoquet, C.; Mourad, R. Probabilistic Graphical Models for Next-generation Genomics and Genetics. In *Probabilistic Graphical Models for Genetics, Genomics, and Postgenomics*; Oxford University Press, 2014; pp. 1–16. doi:10.1093/acprof:oso/9780198709022.001.0001.
66. Heckerman, D.; Breese, J. Causal independence for probability assessment and inference using Bayesian networks. *IEEE Transactions on Systems, Man, and Cybernetics - Part A: Systems and Humans* **1996**, *26*. doi:10.1109/3468.541341.
67. Borsuk, M.E.; Stow, C.A.; Reckhow, K.H. A Bayesian network of eutrophication models for synthesis, prediction, and uncertainty analysis. *Ecological Modelling* **2004**, *173*. doi:10.1016/j.ecolmodel.2003.08.020.
68. Sevinc, V.; Kucuk, O.; Goltas, M. A Bayesian network model for prediction and analysis of possible forest fire causes. *Forest Ecology and Management* **2020**, *457*. doi:10.1016/j.foreco.2019.117723.
69. E. Neapolitan, R. *Learning Bayesian networks*; Prentice Hall, 2004; p. 433.
70. Kabli, R.; Herrmann, F.; McCall, J. A chain-model genetic algorithm for Bayesian network structure learning. Proceedings of the 9th annual conference on Genetic and evolutionary computation - GECCO '07; ACM Press: New York, New York, USA, 2007. doi:10.1145/1276958.1277200.
71. Scanagatta, M.; Salmerón, A.; Stella, F. A survey on Bayesian network structure learning from data. *Progress in Artificial Intelligence* **2019**, *8*. doi:10.1007/s13748-019-00194-y.
72. L. Zhang, N. COMP538: Introduction to Bayesian Networks Lecture 6: Parameter Learning in Bayesian Networks, 2008.
73. Spiegelhalter, D. Lecture 6: Bayesian estimation, 2016.
74. Fan, Z.; Chin, A. Lecture 20 — Bayesian analysis, 2016.
75. Orlof, J.; Bloom, J. Comparison of frequentist and Bayesian inference, 2014.
76. J. Storkey, A. Machine Learning and Pattern Recognition: Note on Dirichlet Multinomial, 2020.
77. Liu, H.; Wasserman, L. Bayesian Inference. In *Statistical Machine Learning*; Carnegie Mellon University, 2014; pp. 299–305.
78. Alvares, D.; Armero, C.; Forte, A. What Does Objective Mean in a Dirichlet-multinomial Process? *International Statistical Review* **2018**, *86*. doi:10.1111/insr.12231.
79. Kelly, D.; Atwood, C. Finding a minimally informative Dirichlet prior distribution using least squares. *Reliability Engineering & System Safety* **2011**, *96*. doi:10.1016/j.res.2010.11.008.
80. Robert, C.P. Bayesian computational tools. *Annual Review of Statistics and Its Application* **2014**, *1*, 153–177. doi:10.1146/annurev-statistics-022513-115543.
81. Koller, D.; Friedman, F. Bayesian Parameter Estimation. In *Probabilistic Graphical Models*; The MIT Press, 2009; pp. 738–739.

82. Bielza, C.; Larrañaga, P. Bayesian networks in neuroscience: a survey. *Frontiers in Computational Neuroscience* **2014**, *8*. doi:10.3389/fncom.2014.00131.
83. Shimony, S.E. Finding MAPs for belief networks is NP-hard. *Artificial Intelligence* **1994**, *68*, 399–410. doi:10.1016/0004-3702(94)90072-8.
84. Pearl, J. *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*, 1 ed.; Morgan Kaufmann Publishers, INC: San Francisco, 1988.
85. Lozano-Pérez, T.; Kaelbling, K. 6.825 Techniques in Artificial Intelligence (SMA 5504), 2002.
86. Guo, H.; Hsu, W. A Survey of Algorithms for Real-Time Bayesian Network Inference. Technical report, Association for the Advancement of Artificial Intelligence, 2002.
87. Shwe, M.; Cooper, G. An empirical analysis of likelihood-weighting simulation on a large, multiply connected medical belief network. *Computers and Biomedical Research* **1991**, *24*, 453–475. doi:10.1016/0010-4809(91)90020-W.
88. Russell, S.; Norvig Peter. *Artificial Intelligence: A Modern Approach*, 3 ed.; Prentice Hall, 2010; pp. 533–535.
89. National Library of Medicine. National Center for Biotechnology Information, 1988.
90. Edgar, R. Gene Expression Omnibus: NCBI gene expression and hybridization array data repository. *Nucleic Acids Research* **2002**, *30*. doi:10.1093/nar/30.1.207.
91. Barrett, T.; Wilhite, S.E.; Ledoux, P.; Evangelista, C.; Kim, I.F.; Tomashevsky, M.; Marshall, K.A.; Phillippy, K.H.; Sherman, P.M.; Holko, M.; Yefanov, A.; Lee, H.; Zhang, N.; Robertson, C.L.; Serova, N.; Davis, S.; Soboleva, A. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Research* **2012**, *41*. doi:10.1093/nar/gks1193.
92. Love, M.I.; Huber, W.; Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biology* **2014**, *15*. doi:10.1186/s13059-014-0550-8.
93. Varet, H.; Brillet-Guéguen, L.; Coppée, J.Y.; Dillies, M.A. SARTools: A DESeq2- and EdgeR-Based R Pipeline for Comprehensive Differential Analysis of RNA-Seq Data. *PLOS ONE* **2016**, *11*. doi:10.1371/journal.pone.0157022.
94. Conesa, A.; Madrigal, P.; Tarazona, S.; Gomez-Cabrero, D.; Cervera, A.; McPherson, A.; Szczesniak, M.W.; Gaffney, D.J.; Elo, L.L.; Zhang, X.; Mortazavi, A. A survey of best practices for RNA-seq data analysis. *Genome Biology* **2016**, *17*. doi:10.1186/s13059-016-0881-8.
95. Wen, G. A Simple Process of RNA-Sequence Analyses by Hisat2, Htseq and DESeq2. Proceedings of the 2017 International Conference on Biomedical Engineering and Bioinformatics - ICBE 2017; ACM Press: New York, New York, USA, 2017. doi:10.1145/3143344.3143354.
96. Jeong, H.H.; Liu, Z. Are HHV-6A and HHV-7 Really More Abundant in Alzheimer's Disease? *Neuron* **2019**, *104*. doi:10.1016/j.neuron.2019.11.009.
97. Nagarajan, R.; Scutari, M.; Lèbre, S. *Bayesian Networks in R*; Springer New York, 2013. doi:10.1007/978-1-4614-6446-4.
98. Scutari, M. Learning Bayesian Networks with the bnlearn R Package. *Journal of Statistical Software* **2010**, *35*. doi:10.18637/jss.v035.i03.
99. R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017.