

マルコフ連鎖分割アルゴリズムの開発と 社会ネットワークシステムへの応用

名古屋商科大学 経営学部

教授 笹沼 克信

教授 韓 尚憲

1. はじめに

マルコフ連鎖 (Markov chain、または MC) は現実の社会の問題を分析するツールとして幅広く利用されている。恐らく最も有名な適用例は、Google の検索エンジンとなっている PageRank アルゴリズムである。このアルゴリズムでは、各 Website のハイパーリンクを介してのつながり (ネットワーク) をマルコフ連鎖 (MC) システムとみなし、これを解いて定常状態を求め、その定常状態確率に基づいて検索結果を表示する順番を決定している。ソーシャルネットワークもマルコフ連鎖システムとみなすことが可能である。マルコフ連鎖ネットワークを構成する各ノードは、個人のアカウント、それらを結びつけるリンクはアカウント同士の結びつき (例えば交流の頻度) を表すことになる。ソーシャルネットワークをマルコフ連鎖システムとして分析することによって、さまざまな情報や考え、行動が誰を中心としてどのような形で社会に広がっていくかを可視化することができるようになる。そしてマルコフ連鎖の様々な応用の中で近年最も注目されているのは、隠れマルコフモデル (Hidden Markov Model、もしくは HMM) である。HMM を用いると (隠れたマルコフ連鎖を仮定することによって)、観測結果に基づいて、外から見ることでできないシステム内部の (つまり隠れた) 状態を知ることが可能になる。HMM は音声認識やコンピュータによる自動翻訳、隠れた情報の抽出などに用いられており、機械学習における基本技術の一つとなっている。

本研究では、マルコフ連鎖システムを応用した HMM によるアルゴリズムを用いて、現実の社会ネットワークシステムにおいて情報の信頼性 (もしくは虚偽性) がどのように変化していくかについて分析する手法を確立した。HMM を用いた分析は、通常、数値的解析に限定されるが、我々はマルコフ連鎖分割法 (Sasanuma *et al.* [1]) という新しいフレームワークを用いることにより、解析的な分析を行うことを可能にした。これにより、社会システムの中で情報の信頼性・虚偽性がどのように変化するかについて、本質的な特徴・性質を知るができるようになった。このような解析的アプローチによる分析は我々が知る限りこれまで行われてこなかった。本研究で提案するマルコフ連鎖分割法を用いた HMM の特性を求める手法は、社会システムの中の情報の信頼性の分析のみならず、医学、AI、情報など多岐に渡る分野での応用が期待される。

2. 社会ネットワークシステムにおいて情報を分析する上での課題

ネットワークシステム上には様々な情報が溢れている。その中には重要な情報もあればあまり重要でないものもある。例えばある情報に関連した Website を見つけたい場

合、やみくもに Website を開いても必要な情報はなかなか得られない。そこで、我々は Google の検索エンジンにアクセスし、Google の開発した PageRank アルゴリズムによってランク付けされた Website のリストを閲覧し、リストの上から（即ち重要度の高い順番に）Website を開いて情報を入手しようとする。このように Website をランク付けする検索エンジンは、ネットワークシステム上では必要不可欠なツールであることは言うまでもない。しかし情報の重要度と同程度に、もしくはそれ以上大切なものは、情報の信頼性（もしくは虚偽性）である。例えばソーシャルネットワークにおいて様々な情報が拡散しているが、それぞれについてどの程度信頼することができるのか。信頼性を定量的に求めることは容易ではない。もし PageRank アルゴリズムで行っているのと同様に、信頼性や虚偽性の確率を求めて情報をランク付けすることができたら、信頼性の高いものを優先的に入手することが可能となり、逆に虚偽性の高いものについてはソーシャルネットワーク管理者に情報の削除要請をすることが可能となる。

ソーシャルネットワークに流れる情報の信頼性を定量的に判断するアルゴリズムを構築することは可能だろうか。本研究では情報の信頼性（もしくは虚偽性）を次のように定義する：情報が事実に基づけば信頼性が高く（虚偽性が低く）、虚偽に基づいた情報は信頼性が低い（虚偽性が高い）。ここで、課題となるのが、情報を聞いて事実に基づいているかどうかを判断する（Test と呼ぶ）ことは可能だが、実際に情報が事実に基づいているかそれとも虚偽に基づいているかは容易に知ることができない、という問題である。なぜなら、情報の発信者は、その情報が事実に基づいているか虚偽に基づいているかを教えてくれないからである。情報が事実に基づいているか、虚偽に基づいているかを知る方法として、隠れマルコフモデル（HMM: Hidden Markov Model）がある。HMM では直接観測できない状態（ここでは、真実と虚偽という状態）が遷移し、マルコフ連鎖を形成している時に、これらの状態からの出力（ここでは、情報）を観測することによって観測できない状態を調べる数学的なツールである。

HMM については、EM アルゴリズムによる数値的な分析（Baum *et al.* [2] 及び Dempster *et al.* [3] を参照）を行うことが可能であり、AI をはじめとして多方面で HMM を用いた分析が行われている。しかし、これらの数値解析的アプローチは計算量の負荷が大きい。ソーシャルネットワーク上に流れる情報の一つ一つに HMM を適用して数値的な分析を行い、それぞれの信頼性を見積もることは不可能ではないかもしれないが、現実的とは言えない。もし解析的に HMM による分析を行い、情報の信頼性もしくは虚偽性を見積もることができれば、数値解析的アプローチと比較してはるかに効率良く信頼性・虚偽性を見積もることが可能となると考えられるが、隠れマルコフモデル（HMM）に対する解析的な分析については、我々が知る限りほとんど存在しない。

HMM の解析的な分析を難しくしている最大の要因は、出力（Emission）の存在である。HMM によるマルコフ連鎖は出力を伴っているために、HMM 全体として見た時に、通常のエルゴード的なマルコフ連鎖の形をしていない。さらに、出力を状態としてマルコフ連鎖のシステムの一部だとみなそうとすると、出力の示す状態はマルコフ連鎖のどの状態に起因して生じたものかが区別されていないという問題にぶつかる。この二つの問題は、

HMM を通常のエルゴード的なマルコフ連鎖と異なるものとさせており、その結果、HMM を解析に解くことを困難にしている。

我々は、以下のような工夫して、HMM 中に存在する出力をマルコフ連鎖に組み込み、HMM のシステム全体を通常のエルゴード的なマルコフ連鎖に変換する。まず、1) 出力も HMM の一部とみなし、出力が有限の時間起きていると考える。(出力の時間が無限小の極限で、従来の HMM に一致する。) 次に、2) 出力を状態とみなし、その起源によって区別した。すなわち、情報を観測した時にその情報を虚偽だと考えることを Test Positive と呼ぶことにすると、Test Positive の確率 (情報を虚偽だと考える確率) は、True Positive (情報は虚偽に基づいている場合) と False Positive (情報は事実に基づいている場合) という二つの状態の和であると考えた。つまり Test Positive という判断をした場合に、True Positive の確率が大きければ、情報は虚偽である可能性が高いことになる。逆に False Positive の確率が大きければ、虚偽だと考えていた情報が実は事実に基づいている可能性が高い、ということになる。以上、1 と 2 の二点の工夫によって、観測不可能な隠れたマルコフ連鎖と観測可能な出力という異なるシステムを組み合わせた従来の HMM を、観測不可能な部分と観測可能な部分の二つが一体化した (通常のエルゴード的な) マルコフ連鎖に変換することができた。このマルコフ連鎖について、マルコフ連鎖分割法 (Sasanuma *et al.* [1]) という新しいフレームワークを用いて解析に解き、True Positive の確率と False Positive の確率を導いた。これにより、情報の虚偽性 (もしくは信頼性) を定量的に分析することが可能となった。

3. 出力を組み込んだ連続時間隠れマルコフモデル (CT-HMM) の設計

本研究では、従来の隠れマルコフモデルを、出力を組み込んだ連続時間マルコフ連鎖 (CTMC: Continuous Time Markov Chain) を用いた隠れマルコフモデル (CT-HMM: Continuous-Time Hidden Markov Model) に拡張し、解析的な分析を行った。CT-HMM を用いることにより、1) 状態遷移が連続時間的に行われる場合の分析が容易となり、また、2) CTMC に対してモデルパラメタを設定すれば、パラメタに依存した観測結果、及び CT-HMM の特性を解析的に分析可能となる。本研究では、情報の信頼性と虚偽性を調べるために、次のような CTMC によるモデルを設計し、分析手法の効果についての検証を行った。

出力を組み込んだ連続時間隠れマルコフモデル (以下、CT-HMM と記述する) は図 1 のように構成される。各状態を次のように定義する (FP, TP, TN, FN は出力を表す状態となっている。) Neg:Negative; Pos:Positive; TP:True Positive; FP: False Positive; TN:True Negative; FN:False Negative

図 1 の p と q は Neg (事実が存在している状態) と Pos (虚偽が存在している状態) の間のポアソン過程の遷移レート (単位時間当たりの平均的な遷移回数) を表す。 γ と μ はそれぞれ出力の状態と出力の終了した状態への遷移のレートを表す。(なお、レートの逆数は、各状態に継続して存在する時間の平均値となる。) α と β は情報が事実に基づい

ているか虚偽に基づいているかを判断する Test（例えば機械学習による予測アルゴリズムを用いた Test；外生的に与える）の誤りの確率を示しており、それぞれ false positive の確率（偽陽性の確率、Type I エラー）と false negative の確率（偽陰性の確率、Type II エラー）を示している。

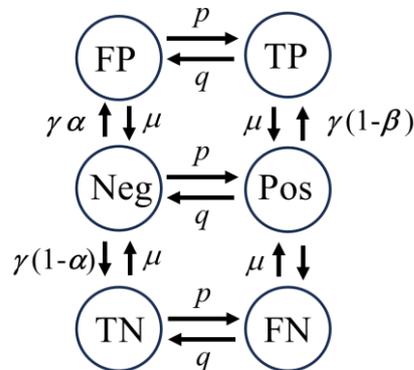


図 1. 出力を組み込んだ連続時間隠れマルコフモデル (CT-HMM)

図 1 で示されたモデル (CT-HMM と呼ぶ) では、Negative (事実が存在している状態) と Positive (虚偽が存在している状態) の二つの (観測不可能な状態) 間を状態が遷移し、その状態を調べるために Test を行い、その結果 (すなわち情報が虚偽か事実に基づいているかという判断) を観測する、というような状況を想定している。Test を行うという行為は、従来の HMM における出力に対応している。Test の結果は、Test Negative と Test Positive の二つのどちらか一方であり、これらの結果は観測可能であるとする。ただし、真の状態が Negative (事実) であったとしても、Test の結果は Test Negative (情報は事実に基づいていると判断される) とは限らず、Test Positive (情報は虚偽に基づいていると判断される) ということもある。同様に、真の状態が Positive であったとしても、Test の結果は Positive であるとは限らず、Negative である場合もある。

このように、Test を行った結果、Test Negative (TestN:これには TN と FN の両方が含まれ、情報を観測する時にはこれらの区別ができない) か、Test Positive (TestP:これには TP と FP が含まれているが、情報を観測する時にはこれらの区別ができない) のどちらかを得ることになる (図 2 を参照)。

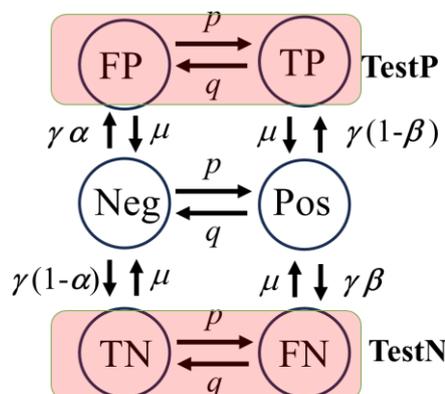


図 2. Test Positive (TestP) と Test Negative (TestN) に対応する状態

一方で、真に Positive (RealP:Real Positive と呼ぶ;虚偽が存在する状態) であるか、真に Negative (RealN:Real Negative と呼ぶ;事実が存在する状態) であるかについては、直接観測できない。例えば FP と TP は観測結果から RealN に属しているのか RealP に属しているのか区別ができないので、観測できない (図 3 を参照)。

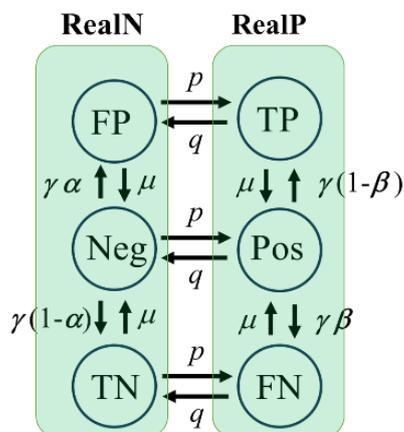


図 3. Real Positive (RealP) と Real Negative (RealN) に対応する状態

我々の CT-HMM (図 1) は通常のエルゴード的マルコフ連鎖であるため、解析的に解くことができる。また、ある状態が実現した後、次にどのような状態が実現するかという情報を簡単に得ることができる。これらの情報を用いて、直接観測することのできる Test 結果 (TestP と TestN) から、観測できない真の状態 (RealP と RealN) の情報を得ることができる。すなわち、情報を観測することによって、観測できない状態、すなわち情報が事実に基づくものか、もしくは虚偽に基づくものかを定量的に (確率的に) 判断することが可能となる。

4. CT-HMM の解析

CT-HMM (図 1) は通常のエルゴード的 CTMC であり、これを解析的に解くために、Sasanuma *et al.* [1] によるマルコフ連鎖分割法を利用した。分析項目は、大別して二種類に分けられる。

第一の種類は、定常状態の特性に関する分析である。例えば TestP、TestN、TP、FP、TN、FN などの状態を取る確率が、偽陽性の確率 α 、偽陰性の確率 β 、そして Test 時間などのモデルパラメタに依存してどのように変化するかを知ることができる。表 1 は Test Negative という結果が得られた場合に、その結果が正しい (つまり True Negative である) 確率を示している。モデルパラメタは、 $p = 0.1$ 、 $q = 1$ 、 $\gamma = 1/12$ 、 $\mu = 2$ とした。すなわち、Negative に平均 10 時間存在した後、Positive に遷移してその状態に平均 1 時間存在し、また Negative に戻ることを想定し、Test は約 12 時間ごとに行われて約 30 分かかるものとしている。 α や β の増大によって TN の割合が低下することが表 1 の結果から確認できる。すなわち、Test (情報の判断) にかかわる誤差が増大すると、Test Negative という判断をした場合において、True Negative (情報が事実に基づく) の可能性が下がって行くことが分かる。

TN/(TN+FN)	alpha											
	0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1	
beta	0	0.968	0.968	0.968	0.968	0.968	0.968	0.968	0.968	0.968	0.968	0.897
	0.1	0.961	0.961	0.960	0.959	0.957	0.955	0.952	0.947	0.937	0.909	0.323
	0.2	0.955	0.954	0.952	0.950	0.947	0.943	0.937	0.927	0.909	0.860	0.323
	0.3	0.949	0.947	0.944	0.941	0.937	0.931	0.923	0.909	0.884	0.819	0.323
	0.4	0.943	0.940	0.937	0.933	0.927	0.920	0.909	0.892	0.860	0.783	0.323
	0.5	0.937	0.934	0.930	0.925	0.918	0.909	0.896	0.876	0.839	0.753	0.323
	0.6	0.931	0.927	0.923	0.917	0.909	0.899	0.884	0.860	0.819	0.726	0.323
	0.7	0.926	0.921	0.916	0.909	0.900	0.889	0.872	0.846	0.800	0.702	0.323
	0.8	0.920	0.915	0.909	0.902	0.892	0.879	0.860	0.832	0.783	0.681	0.323
	0.9	0.914	0.909	0.903	0.894	0.884	0.869	0.849	0.819	0.768	0.662	0.323
	1	0.909	0.903	0.896	0.887	0.876	0.860	0.839	0.806	0.753	0.645	0.323

表 1. Test Negative (TestN) の場合に True Negative (TN) である確率

第二の種類は、状態の遷移に関する分析である。CT-HMM は通常の連続時間マルコフ連鎖なので、どの状態からどの状態に遷移するか、(例えば CT-HMM に対応する Embedded 離散時間マルコフ連鎖を用いて) その確率を解析的に計算することが可能である。例えば、True Negative (TN) という状態が実現した場合に、次に False Negative (FN) の状態が続く確率を解析的に求めることができる。このようにして、(TN, TN)、(FN, FN)、(TN, FN)、(FN, TN) の 4 つの場合についてそれぞれの組み合わせ (事象系列と呼ぶ) が起こる確率を解析的に求めることが可能となる。これら 4 つの組み合わせは、観測可能な Test 結果で言えばどれも Test Negative が連続して起こる事象系列に対応しているが、真実の状態の事象系列は全て異なっている (例えば、(TN, FN) は (Neg, Pos) が真の状態の事象系列であるが、(FN, TN) は (Pos, Neg) が真の状態の事象系列となる)。さらにこれらの 4 つの結果から、Test Negative (TestN ; TN か FN のどちらかの結果が得られる状態) が連続する事象系列の確率や Test Positive (TestP ; TP か FP のどちらかの結果が得られる状態) が連続する事象の確率も解析的に求められる。

表 2 は TestP を観測した場合にその次も TestP が起こる条件付確率が、状態間の遷移レート p と q にどのように依存しているかを示したものである。(ここでは、 $\alpha=0.1$ 、 $\beta=0.3$ 、 $\gamma=1/12$ 、 $\mu=2$ としている。) p の増加に伴い条件付確率が上昇し、 q の増加に伴い条件付確率が減少する様子を観察できるが、これは p/q 比の増加に伴い真の状態が Positive である可能性が増えるためであり、結果は直感的に明らかである。ここで注目したいのは、 p/q 比が一定であっても p と q の大きさに依存して条件付確率が異なる、という点である。そして、その依存関係は、我々のアプローチによって解析的に理解することが可能となった。例えば、 p と q が大きい極限では、条件付確率は p/q 比のみに依存することが解析的に示される。

Prob(TestP TestP)	p											
	0.01	0.03	0.1	0.3	1	3	10	30	100	300	1000	
q	0.01	0.58	0.63	0.66	0.68	0.69	0.70	0.70	0.70	0.70	0.70	0.70
	0.03	0.43	0.53	0.60	0.65	0.68	0.69	0.70	0.70	0.70	0.70	0.70
	0.1	0.23	0.34	0.46	0.57	0.65	0.68	0.69	0.70	0.70	0.70	0.70
	0.3	0.14	0.19	0.29	0.42	0.57	0.65	0.68	0.69	0.70	0.70	0.70
	1	0.11	0.12	0.16	0.25	0.40	0.55	0.65	0.68	0.69	0.70	0.70
	3	0.10	0.11	0.12	0.16	0.25	0.40	0.56	0.65	0.68	0.69	0.70
	10	0.10	0.10	0.11	0.12	0.15	0.24	0.40	0.55	0.65	0.68	0.69
	30	0.10	0.10	0.10	0.11	0.12	0.15	0.25	0.40	0.56	0.65	0.68
	100	0.10	0.10	0.10	0.10	0.11	0.12	0.15	0.24	0.40	0.55	0.65
	300	0.10	0.10	0.10	0.10	0.10	0.11	0.12	0.15	0.25	0.40	0.56
	1000	0.10	0.10	0.10	0.10	0.10	0.10	0.11	0.12	0.15	0.24	0.40

表 2. Test Positive (TestP) の場合に、次も TestP が続く確率

隠れマルコフモデルの分析は従来、アルゴリズムによる計算科学的なアプローチが取られていた。しかし、表1や表2で例示したように、本研究による解析的な分析により、事象系列が起きる確率がモデルパラメタによってどのように変化するか知ることが可能となった。観測された Test の事象系列から最も有り得る可能性の高い隠れた状態の並びを求める場合、従来は Viterbi アルゴリズム [4] を用いて数値的に探索して最適解を求めることに留まるが、我々の場合はモデルパラメタに依存して最も確からしい隠れた状態の並びがどのように変化するか知ることができる。また、観測された Test の事象系列から最も有り得る可能性の高いモデルパラメタ（状態間の遷移確率と出力確率）の予測値を求める場合、従来は Baum-Welch アルゴリズム [2] を用いて数値的にパラメタを探索するが、我々の場合は事象系列が起こる確率の解析解が得られているので、モデルパラメタを解析的に（つまり、方程式を解くことによって）求めることが可能となる。このような試みは、我々が知る限り従来無かったものである。

5. 情報の虚偽性・信頼性のパラメタ依存性

虚偽（偽りの事実）は長い時間存在しているわけではなく、平均的には短い時間で（虚偽だということが知られて）価値を失い消滅すると考えられる。ここでは、虚偽の存続する時間を変えて情報の虚偽性（どの程度情報が虚偽に基づいているか、という確率）を求めた。具体的には、 $p=0.1$ 、 $\gamma=1/12$ 、 $\mu=2$ 、 $\alpha=0.1$ 、 $\beta=0.3$ 、とし、 q だけ変更を加えて虚偽（偽りの事実）の存続する時間を変えて、Test Positive が連続した場合の True Positive（虚偽性）と False Positive（信頼性）を評価した。虚偽の存続する時間が長いほど、そして Test Positive の回数が多いほど、True Positive と判断しやすくなることがわかる。

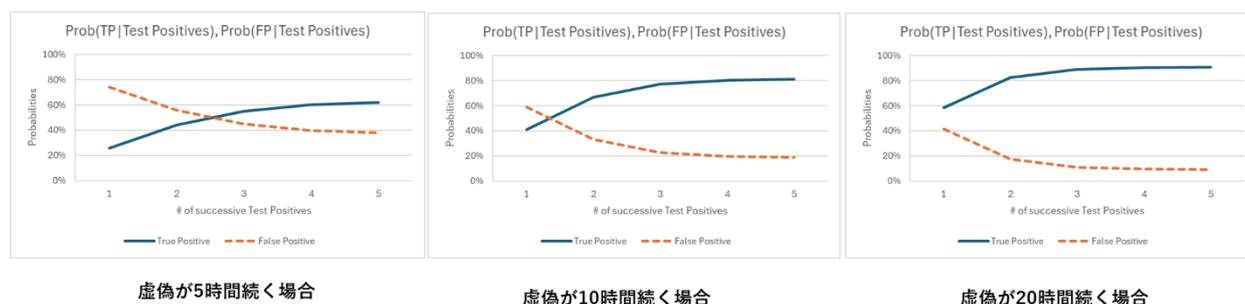


図 4. True Positive（情報が虚偽に基づく確率）と False Positive（情報が事実に基づく確率）

5. まとめと今後の課題

ソーシャルメディアにおいて拡散される情報は、事実に基づくものと、虚偽に基づくものが混在している。本研究では、情報の信頼性を定量的に分析するために、従来の隠れマルコフモデル (HMM) を新しい連続時間隠れマルコフモデル (CT-HMM) に拡張し、Sasanuma *et al.* [1]によるマルコフ連鎖分割法を用いて解析を行った。これにより、情報の信頼性、または虚偽性を効率良く調べることが可能となった。HMM によって分析を行う場合は、従来は数値的なアルゴリズムによるアプローチを取ることが一般的であるが、マルコフ連鎖分割法を用いることにより解析的に分析することが可能となった。

本研究で用いた CT-HMM は観測不可能な状態数が 2 つ、すなわち事実 (Negative) と虚偽 (Positive)、そして観測結果が 2 つ、すなわち事実に基づくと判断される情報 (Test Negative) と虚偽に基づくと判断される情報 (Test Positive) というものであったが、これをそれぞれ 3 つ以上に拡張することは容易である。今後は CT-HMM を拡張して実際のソーシャルメディア上の情報について、分析を進める。また、本研究で開発した手法を医学、AI、情報など様々な分野に応用が可能である。本研究の成果は論文の形でまとめられており、出版される予定である。

謝辞

本研究を遂行するにあたり、公益財団法人 天野工業技術研究所から多大なご支援を頂きました。ここに記して謝意を示します。また、ニューヨーク州立大学ストーニーブルック校医学部教授 James Dilger 博士、そして名古屋商科大学経営学部の濱砂稔真氏に協力を頂きました。ここに記し、深く感謝申し上げます。

参考文献

- 1) Sasanuma, K., Hampshire, R., and Scheller-Wolf, A., Markov chain decomposition based on total expectation theorem. <https://arxiv.org/abs/1901.06780>
- 2) Baum, L.E., Petrie, T., Soules, G. and Weiss, N., 1970. A maximization technique occurring in the statistical analysis of probabilistic functions of Markov chains. *The annals of mathematical statistics*, 41(1), pp.164-171.
- 3) Dempster, A.P., Laird, N.M. and Rubin, D.B., 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the royal statistical society: series B (methodological)*, 39(1), pp.1-22.
- 4) Viterbi, A., 1967. Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE transactions on Information Theory*, 13(2), pp.260-269.