

物語イベントの状態遷移解析：半教師あり GMM を用いた 転調点抽出

——宮沢賢治『銀河鉄道の夜』を例に

要旨

本研究は、物語を「状態」と「遷移」として定式化し、転調点を再現可能な量として抽出する枠組みを提示する。宮沢賢治『銀河鉄道の夜』のイベント列を、心理特徴量 $X=(m, iso)$ へ写像し、少量ラベル（アンカー）を用いた半教師あり GMM を EM で推定する。各イベントの事後確率エントロピーを boundaryness と定義し、上位点を転調点候補として抽出する。モデル選択（BIC/AIC/Silhouette）と alpha 比較（アンカー整合、境界安定、遷移数、entropy 総量）により、二状態モデルが解釈性と指標の両面で最適であることを示す。さらに、転調点は遷移回数と一致せず、心理的不確定性の極大として分離して現れることを明らかにする。

キーワード：半教師あり学習、ガウス混合モデル、EM アルゴリズム、転調点、エントロピー、デジタル・ヒューマニティーズ

1. 背景

物語を「状態」と「遷移」として定式化する試みは、解釈の再現可能性を担保するデジタル・ヒューマニティーズの基盤である。しかし現状の物語解析は、(i) 章・場面の区切りを人手に依存する、(ii) 分割点の意味づけが事後的で恣意的になりやすい、(iii) 教師データ不足により監督学習が成立しにくい、という制約を併せ持つ。この制約は、局面の転調が重要である一方でラベル付けが高コストであり、かつ解釈可能性が要求される文学作品で顕在化する。ゆえに、少量のラベルで状態の意味を固定しつつ、転調点をモデル内部量として抽出できる手法が必要である。（Blei, 2012）

2. 手法

2.1 データ（イベント列）

入力 は 章（scene_id）ごとに分割されたイベント列であり、各イベントに根拠引用（evidence）と要約（desc）を付与する。イベントは全体時系列（global_step）として連番管理し、章開始位置は時系列上の参照線として利用する。

2.2 特徴量設計

各イベントを、以下の 2 次元心理特徴量で表現する。

- m (Morality) : 自己防衛・自己中心への退避 (0) から、共感・他者配慮・共同幸福 (1) への連続値。
- iso (Isolation/Survival Urgency) : 安心・切迫感の低さ (0) から、孤独・欠乏・喪失・切迫 (1) への連続値。
- Time (参照軸) : global_step を正規化して 0-1 に写像し、可視化の参照軸として用いる。

2.3 半教師あり GMM (EM)

観測 $X=(m, iso)$ に対し、 K 成分の GMM を EM で推定する。少数の labeled 点（アンカー）を与え、unlabeled 点と混合して推定する。アンカーの寄与は係数 α で重み付けし、 α が大きいほどクラスラベル（解釈ラベル）が固定化される。推定は複数初期値で行い、seed に対する安定性を評価する。（Bishop, 2006; Dempster et al., 1977; Chapelle et al., 2006）

2.4 転調点 (boundary) の定義

各イベント i について GMM の事後確率 $p(z|x_i)$ を算出し、そのエントロピー H_i を boundaryness と定義する。 H_i が大きい点は状態帰属が曖昧な点であり、転調点候補である。本研究では entropy 上位 Top20 を boundary として抽出し、時系列上に配置する。（Bishop, 2006; Settles, 2009）

$$H_i = - \sum_{k=1..K} p(z=k | x_i) \log p(z=k | x_i)$$

2.5 評価指標

(i) K 選択：BIC/AIC/Silhouette により適合と分離を比較し、目的（転調抽出）に照らして解釈可能性を優先する。(ii) α 効果：アンカー整合、boundary の seed 安定性、 $\alpha=0$ に対する境界重なり、遷移回数、entropy 総量を比較する。

表 1: ルーブリック (m/iso/Time の定義と採点規則)

特徴量	0 の意味	1 の意味
m (Morality)	0 = 自己防衛・自分中心への退避	1 = 共感・他者救済・自己犠牲の受容
iso (Isolation/ Survival_Urgency)	0 = 安定（切迫感が低い）	1 = 孤独・欠乏・喪失が差し迫る（存在的危機）
Time (参考)	0 = 物語の冒頭	1 = 物語の終盤（global_step 正規化）

表 2: K 別モデル選択指標 (BIC/AIC/Silhouette/Anchor agreement)

K	BIC (↓)	AIC (↓)	Silhouette (↑)	Anchor agreement
1.0	-197.66039	-208.00026	nan	0.714286
2.0	-197.834997	-221.099705	0.410386	0.857143
3.0	-182.828083	-219.017627	0.362225	1.0
4.0	-175.279388	-224.39377	0.245848	0.857143
5.0	-156.319921	-218.359141	0.362089	1.0
6.0	-140.541032	-215.505089	0.38262	1.0

表 3: α 別比較 (anchor_acc, boundary_Jaccard, transitions, entropy_sum 等)

alpha	anchor_a cc	boundary _Jaccard	overlap_v s0	transition s	entropy_s um	seeds
0.0	0.857	1.0	1.0	14.0	4.244	5.0
1.0	1.0	0.76	0.495	20.4	16.665	5.0
5.0	1.0	1.0	0.905	12.0	3.694	5.0
10.0	1.0	1.0	0.905	12.0	3.152	5.0
20.0	1.0	1.0	0.818	12.0	3.199	5.0
50.0	1.0	1.0	0.667	16.0	3.47	5.0

図 1: m-iso 散布図 (点サイズ=entropy、円=boundary Top20)

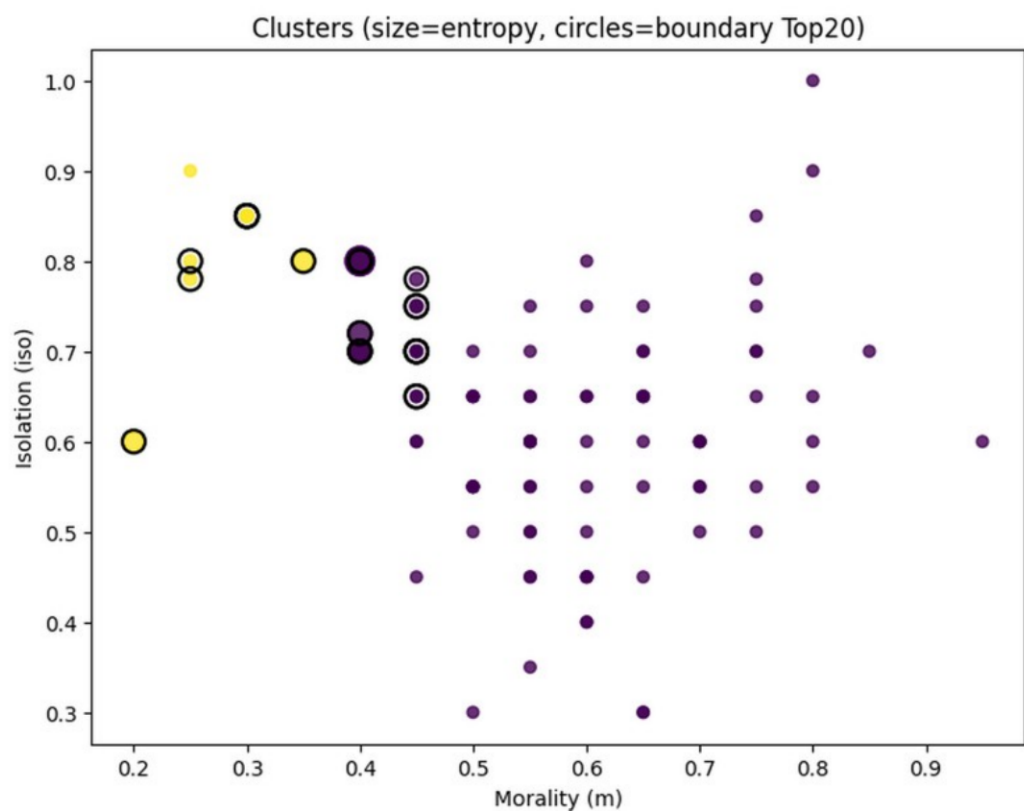


図 2: クラスタ中心と共分散楕円 (2σ)

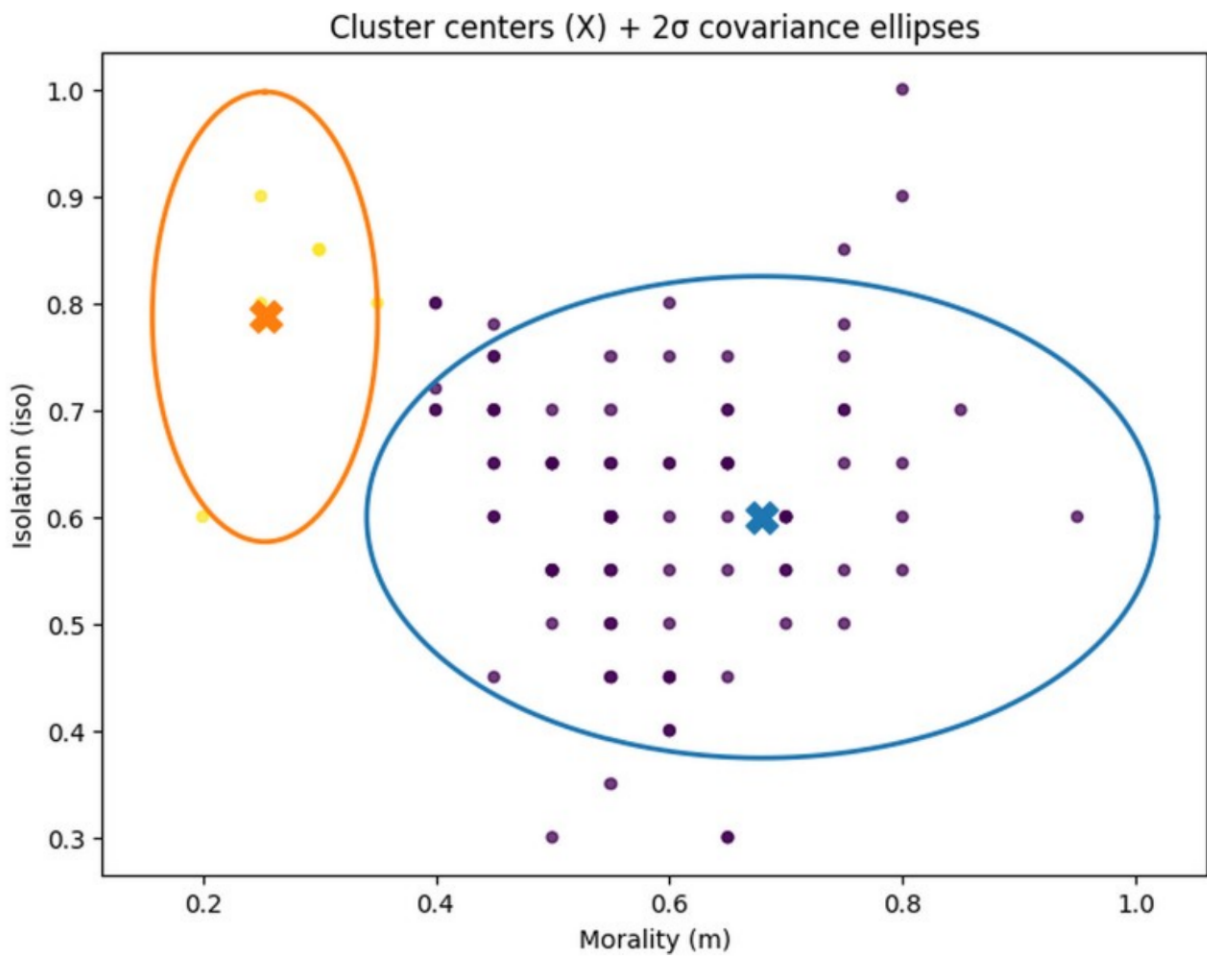


図 3: boundaryness 時系列（縦線=章開始）

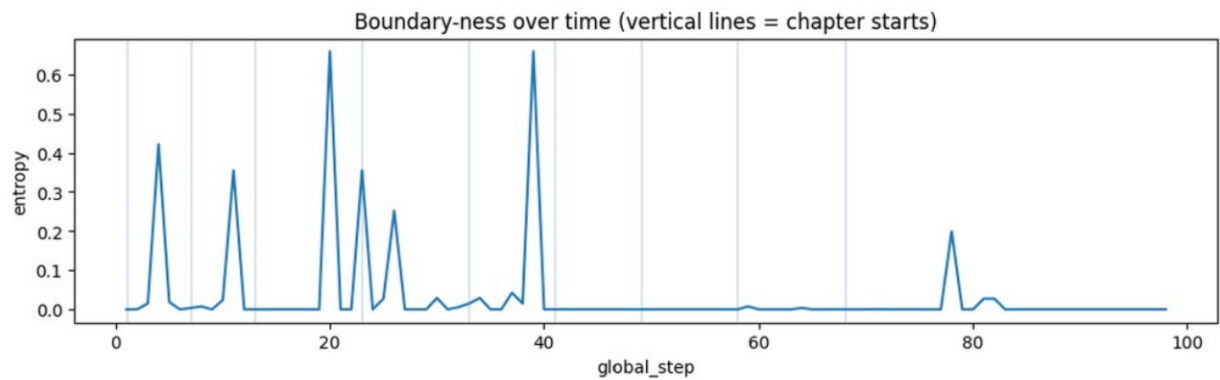


図 4: 状態（クラスタ意味）の時系列

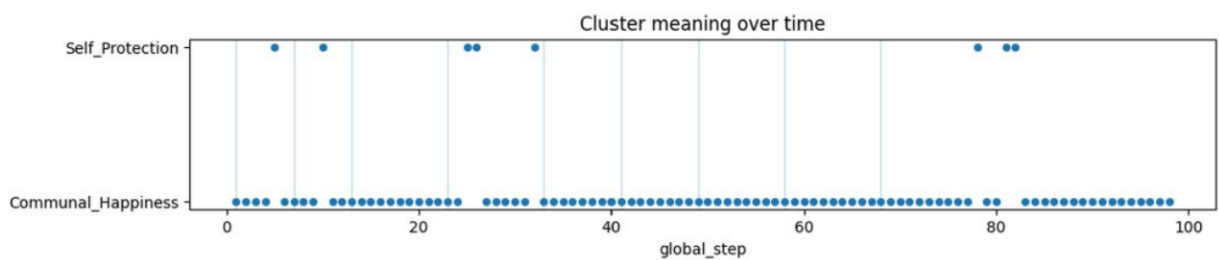


図 5: 章別 turning density (entropy 総量)

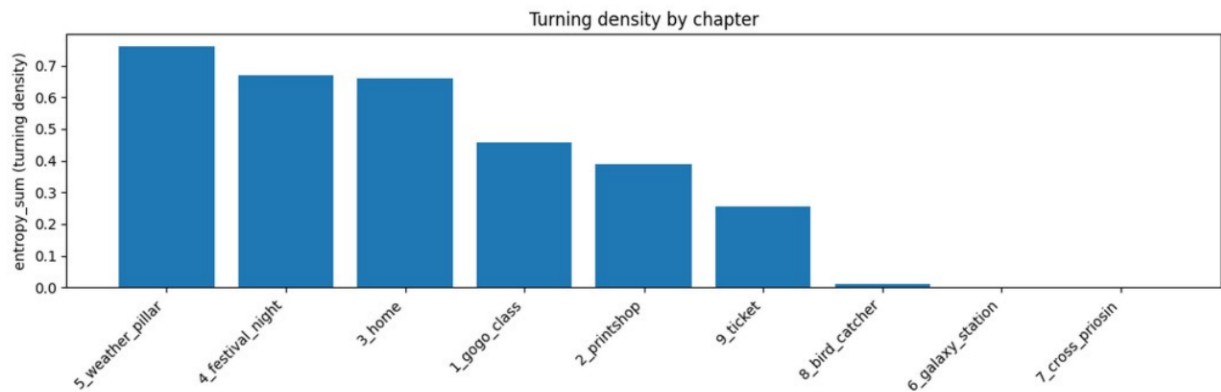
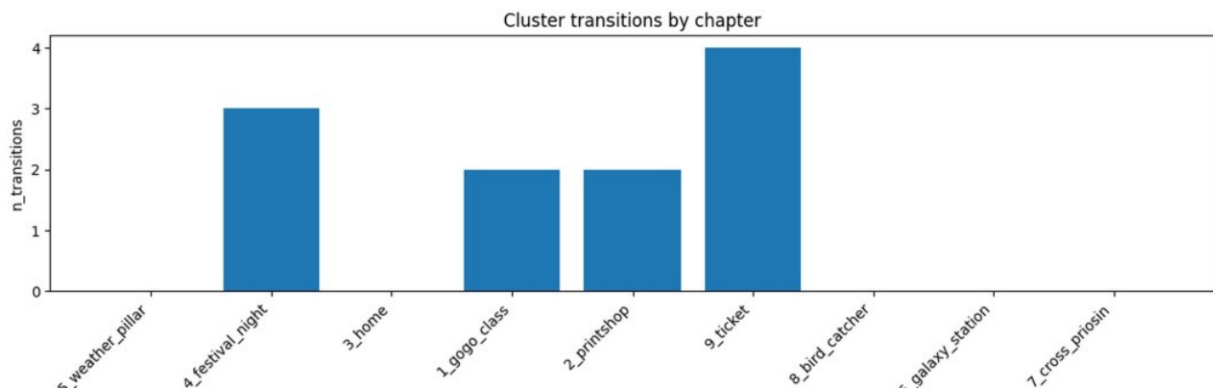


図 6: 章別クラスタ遷移回数



3. 結果と考察

3.1 状態数 K は 2 が最適である

K=1 は状態分割が生じず、遷移・転調の定義が成立しない。K>=3 は過剰分割により転調と状態の微分化が混同される。指標比較では BIC 最小かつ Silhouette 最大が K=2 であり、二状態構造が最も妥当である。

3.2 二状態の意味は「共同幸福」と「自己防衛」に対応する

m/iso 空間では低 m・高 iso 近傍に一群、高 m・中 iso 近傍に一群が形成される。アンカー導入により両群の意味づけが「自己防衛 (Self_Protection)」と「共同幸福 (Communal_Happiness)」へ固定され、状態遷移が解釈可能な系列として読める。状態は章ではなく心理局面に対応し、章内でも局面は揺れ、遷移は局所的に発生する。

3.3 転調点は時系列上でスパイクとして現れる

boundaryness 時系列は複数の明確なスパイクを持つ。これは、どちらの状態にも属し切らない観測が局所的に生じていることを意味する。転調点は遷移回数と必ずしも一致せず、同じ章でも「決断としての明瞭な切替」は entropy を上げにくく、「葛藤・逡巡・曖昧な局面」は entropy を上げる。ゆえに転調点は心理的不確定性の極大として定義される。

3.4 章別集計は「転調密度」と「遷移回数」の違いを可視化する

章別 turning density (entropy 総量) は特定章に偏在する一方、章別遷移回数は別の章で最大化し、両者は一致しない。この差は、(i) 価値判断や切迫感が揺れ続ける高不確定章と、(ii) 状態が明確に切り替わる決断章が分離して存在するという二層性を示唆する。

3.5 alpha は「発見」と「解釈固定」のトレードオフを支配する

alpha=0 は boundary Top20 の seed 間一致が高く転調点の再現性が高い一方、アンカー整合は完全ではない。alpha=5~20 はアンカー整合が高く境界安定性も維持されるが、alpha=0 との重なりは低下し探索的な転調点が抑制される。alpha=1 は遷移回数と entropy 総量が過大化し不安定である。したがって alpha は単なるハイパーパラメータではなく、発見志向と解釈固定志向のいずれを優先するかをコード化する設計変数である。

4. むすび

『銀河鉄道の夜』のイベント列は m/iso の 2 次元空間で二状態として安定に表現でき、事後確率エントロピーの極大点として転調点を定義することで、章構造と整合する形で状態遷移と転調の定量抽出が可能である。

5. 限界と今後の課題

本研究は、ループリック設計に内在する解釈バイアス、単一アノテータに起因する主観性、2 次元圧縮による情報落ちの影響を受ける。今後は多者アノテーションによる一致度評価、特徴量拡張（語り手距離、時間知覚、象徴密度等）、他作品への適用による外的妥当性検証が必要である。

参考文献

本研究の再現・検証に用いた基礎資料（オンライン）は以下のとおりである。

ginga-narrative-gmm（再現用リポジトリ）. <https://github.com/Mokafe/ginga-narrative-gmm>（最終アクセス: 2026-01-02）.

ginga_step_log_script.md（ステップログ可視化スクリプト）. https://github.com/Mokafe/ginga-narrative-gmm/blob/main/docs/ginga_step_log_script.md（最終アクセス: 2026-01-02）.

Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.

Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of the Royal Statistical Society: Series B, 39(1), 1-38.

Chapelle, O., Schölkopf, B., & Zien, A. (Eds.). (2006). Semi-Supervised Learning. MIT Press.

Settles, B. (2009). Active Learning Literature Survey. University of Wisconsin–Madison, Computer Sciences Technical Report 1648.

Blei, D. M. (2012). Topic Modeling and Digital Humanities. *Journal of Digital Humanities*, 2(1).