051

# 統計的手法と機械学習に基づく地すべりおよび土砂崩壊地点の特徴評価

Evaluation of the characteristics at the landslide occurrences based on statistical methods and machine learning

12TM0300 古川 昭太 Shota Furukawa 指導教員 丸山 喜久

## SYNOPSIS

This study evaluated the topographic characteristics at the landslide outbreak locations based on the covariance structure analysis and the machine learnings. The support vector machine (SVM) and random forest (RF) were employed in this study to achieve the objective. The estimated results were compared to investigate the accuracy of the detection of landslide occurrences. The SVM and RF gave reasonable results if a technique for oversampling is applied to correct for a bias in the training dataset. Lastly, the discriminated results were projected onto a map to evaluate the possibility of landslide occurrences.

1. はじめに

わが国では、自然災害の発生に伴って多くの地すべりの 被害を被ってきた.2004年10月23日に発生した新潟県中越 地震では地すべり131ヶ所、崖崩れ115ヵ所、土石流21ヵ所 と多くの土砂災害の被害が出た<sup>1)</sup>.また2011年東北地方太 平洋沖地震では北海道と青森を除く東北地方と群馬・栃 木・茨城・埼玉・群馬で大規模な範囲での斜面崩壊が確認 された<sup>2)</sup>.その他、広島県広島市北部で発生した平成26年8 月豪雨<sup>3)</sup>など、豪雨による土砂災害も多く発生している. 地震調査委員会では今後30年間でマグニチュード7クラス の地震が首都圏で起こる確率は70%としており、大規模な 土砂災害の発生も懸念される.このことから、どのような 地形的特徴が地すべりの発生の原因になり得るのかをあら かじめ予測しておき、地すべりの危険性が高い地域を予測 することは有意義であると考えられる.

そこで本研究では、まず複数の構成概念間の関係を検討 できる統計的手法である共分散構造分析 (Covariance structure analysis)<sup>4)</sup>を行い、どの素因が地すべりや土砂崩壊 の発生に大きく影響するのかを明らかにする. さらに、比 較的少数の説明変数を用いたサポートベクターマシン (Support vector machine)<sup>5)</sup>とランダムフォレスト(Random forest)<sup>6)</sup>に基づく機械学習を実行し、地すべりや土砂崩壊発 生地点を予測するモデルを構築することを目的とする. 分 析手法ごとの結果と既存の地すべり地形分布図と比較し精 度を評価することによって、今後どのような地域で地すべ り発生が起こり得るのかを予測することを目指す.

# 2. 対象地域と使用データ

本研究の対象地域は宮城県仙台市と2016年熊本地震の影響を受けた熊本県とした.使用した主なデータは国土交通 省が整備する国土数値情報<sup>7)</sup>,微地形区分<sup>8)</sup>,防災科研の地 すべり地形分布図<sup>9)</sup>である.

国土数値情報は標高・傾斜角などの地形データが5次メッ

シュ (250 mメッシュ) ごとに格納されているものを使用 した.地すべり地形分布図<sup>9)</sup>は,地すべり地形を航空写真 から判読したものである.宮城県仙台市の地すべり地形分 布を図-1に示す.



図-1 宮城県仙台市の地すべり地形分布図

#### 3. 共分散構造分析に基づく地すべり発生確率の評価

共分散構造分析では、素因(観測変数)が誘因(潜在変数) の影響を受けるものとし、さらに、それらがいくつかの要 因に集約できると仮定する.要因と地すべり発生の確率と の相関を評価する.基本パスモデルを図-2に示す.なお、 基本パスモデルは既往研究<sup>10)</sup>を参考に仮定した.



図-2 本研究で仮定した基本パスモデル

メッシュごとの地すべり発生確率をベイジアンモデルと ファジーセットモデルの二つの手法により求める<sup>11)</sup>. ベイ ジアンモデルの発生確率の計算式を式(1)に示す.

$$S_{q} = \frac{P\{c_{1j}\}...P\{c_{mj}\}}{P\{c_{1j}, c_{2j}, ..., c_{mj}\}} \cdot P\{T_{q}\} \cdot \frac{P\{T_{q}|c_{1j}\}}{P\{T_{q}\}} ... \frac{P\{T_{q}|c_{mj}\}}{P\{T_{q}\}}$$
(1)

$$P\{c_{ij}\} = \frac{N_{ij}}{A} \tag{2}$$

$$P\{c_{1j},...,c_{mj}\} = \frac{N_x}{A}$$
(3)

$$P\{T_q\} = \frac{N_0}{A} \tag{4}$$

$$P\{T_{q} | c_{ij}\} = \frac{N_{d_{ij}}}{N_{ii}}$$
(5)

をそれぞれ代入する必要がある.ここで、Niiは i 番目の素

因におけるカテゴリ jのメッシュ数, Aは全メッシュ数, N<sub>X</sub> はメッシュ qにおける i~m番目の素因に対する全カテゴリ において等しい属性のメッシュ数, N<sub>0</sub>は地すべり地形のメ

ッシュ数, Nduは i 番目の素因におけるカテゴリ j の領域に

属する地すべり地形のメッシュ数である. また,ファジー セットモデルは式(6)に帰着する.

$$S_{q} = \mu_{s} \left\{ q \middle| c_{1j} ..., c_{ij} \right\} = 1 - \prod_{i=1}^{m} \left\{ 1 - \frac{N_{d_{ij}}}{N_{ij}} \right\}$$
(6)

ベイジアンモデル,ファジーセットモデルの両方で地す ベり発生確率を計算し,図-3の基本パスモデルの適合度を 算出する.なお,共分散構造分析に用いる観測変数の値は 式(5)の条件付き確率とし,平均0,標準偏差1に標準化し てある.さらに,基本モデルから観測変数を1または2個 減らしたモデルから,最適なパスモデルを探索した.

それぞれのモデルの地すべり地形の的中率を算出し,最 も的中率の高いモデルを採用することとする.ここでの的 中率とは,実際に地すべり発生の危険がある地点を「危険 あり」,危険がない地点を「危険なし」とそれぞれ正しく判 別したメッシュ数の和を全体のメッシュ数で除したもので ある.具体的には,横軸にパスモデルをもとに算出した標 準化された地すべり発生確率を,縦軸にその度数および累 積頻度を示す.さらに,負極側から地すべり地形の累積頻 度曲線を,正極側から非地すべり地形の累積頻度曲線を描 く(図-4).つまり,地すべり地形をAグループ,非地すべ り地形をBグループとすると,

$$F_{A}(X) = 1 - \int_{-\infty}^{X} f_{A}(x) dx \qquad \int_{-\infty}^{+\infty} f_{A}(x) dx = 1$$
(7)  
$$F_{B}(X) = \int_{-\infty}^{X} f_{B}(x) dx \qquad \int_{-\infty}^{+\infty} f_{B}(x) dx = 1$$
(8)

となり、  $F_A(X), F_B(X)$ の交点が判別の閾値となる.

図-5 に最も的中率の高かったファジーセットモデルに おける基本モデルから起伏量と土地利用を抜いたモデルの パス図と的中率を示す.パスモデルの相関係数から,地す べりには要因1(標高,起伏量,傾斜角,傾斜方向)が大 きく影響していることがわかる.また,要因1内の相関係 数を見ると,標高と傾斜角の相関係数が大きいことがわか る.採用したパスモデルに用いられている説明変数を使用 して,次章の機械学習を行う.



図-5 採用したパスモデル(的中率: 74.00%)

### 4. 機械学習に基づく地すべり発生箇所の評価

本研究では、共分散構造分析に加えて、機械学習に基づく 分析を実施した.ここでは、サポートベクターマシンとラ ンダムフォレストの二手法を利用した.全データの10%を 学習データとして地すべり発生個所の分類器を作成し、学 習に使用しなかったデータをテストデータとして分類器の 精度評価に利用した.

#### (1) サポートベクターマシン

サポートベクターマシン(SVM)は、教師あり学習を用い たパターン認識手法の一つである.基本的には2つのクラ スを識別するための識別機を構成するための学習手法とさ れており、認識性の優れた学習モデルの一つといわれてい る.

SVM による 2 クラス分類では,式(9)で表される最適化 問題を解く.ここで,w は分離超平面の法線ベクトル, は入力ベクトル x を特徴空間 F へ非線形写像する関数 (式(10)), b はスカラー変数を表す<sup>11)</sup>.

$$\min_{\boldsymbol{w},\boldsymbol{b},\boldsymbol{\zeta}} \frac{1}{2} \| \mathbf{w} \|^{2} + C \sum_{i \in [n]} \zeta_{i}$$
s.t.  $y_{i}(\mathbf{w}^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{i}) + \boldsymbol{b}) \ge 1 - \zeta_{i}, i \in [n], \zeta_{i} \ge 0, i \in [n]$ 

$$K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \boldsymbol{\phi}(\mathbf{x}_{i})^{\mathrm{T}} \boldsymbol{\phi}(\mathbf{x}_{j}) \tag{9}$$
(10)

さらに、本研究では非線形写像のための関数に RBF カーネルを用いる.

$$K(\mathbf{x}_{i}, \mathbf{x}_{j}) = \exp(-\gamma \|\mathbf{x}_{i} - \mathbf{x}_{j}\|^{2})$$
(11)

式(11)の $\gamma$ と式(9)の Cは、ハイパーパラメータと呼ばれ分類 結果に大きく影響する<sup>12)</sup>.

本研究ではグリッドサーチ<sup>13</sup>に基づきハイパーパラメ ータを設定する.グリッドサーチとは、 $C \ge \gamma$ の値の範囲 を任意に設定し、交差検証法により、最適な $C \ge \gamma$ の値を 求める手法である.本研究では、全データの10%を学習デ ータ、 $\gamma \ge 10^{-5} \sim 10^{5}$ ,  $C \ge 10^{-2} \sim 10^{2}$ の範囲に設定して、交 差検証法(k = 10)によって最適なパラメータを求めた.

図-6 に分析結果と既存の地すべり地形分布図との比較 を行い,作成した地すべり地形評価図を示す.また,的中 率(地すべりの有無を正しく分類した割合)と,地すべり 地形のメッシュだけで算出した的中率を示す.的中率は共 分散構造分析に比べ向上したが,地すべり発生地点の抽出 精度が悪く,特に実際に地すべりの危険がある地点を危険 なしと判定する危険個所の見落としが多くみられた.

### (2) ランダムフォレスト

ランダムフォレスト(RF)は、複数の決定木を用いる集団 学習法(アンサンブル学習)の一つである。各決定木での予 測結果を多数決することにより、結果の取得を行う。決定 木学習とは、データの種類に応じて決定木を成長、分類さ せていく学習手法であり、式(12)に示す情報利得 IG が最大 となるようにする<sup>14)</sup>.

$$IG(D_{p}, f) = I(D_{p}) - \sum_{j=1}^{m} \frac{N_{j}}{N_{p}} I(D_{j})$$
(12)

ここで、**D**pは学習データ、N はノード, j は注目している

データ, *I* は不純度, *m* は個々の木を分割するノード数で ある.不純度はデータに偏りがあるほど,大きな値になる 指標であり,エントロピー,ジニ係数,分類誤差などが用 いられる.式(13)はエントロピーの式<sup>15)</sup>を表している.

$$I(N) = -\sum p(k|N) \log p(k|N)$$
(13)

ここで**p(k|N)**は,各ノードNでクラスkを取る確率である. 図-6 に地すべり地形評価図および的中率を示す. SVM による分析と同様に,的中率そのものは向上したものの, 地すべり発生地点の抽出の精度が低く,危険個所の見落と

#### (3) 不均衡データの整形

しが多くみられた.

SVM と RF の結果,地すべり発生地点の見落としが多かったことが課題として挙げられた.その原因としては,本研究で取り扱っている地すべり地形分布図のデータが,正例(地すべり発生が認められるメッシュ)に対して負例(地すべり発生が認められないメッシュ)の数が極端に多い不均衡データであるため,負例を過剰に抽出しまったことが考えられる.そこで,本研究ではこのような不均衡データにおける誤分類を解決するための手法として,サンプリング法の一つである SMOTE (Synthetic Minority Over-sampling Technique)<sup>17)</sup>を適用した.

SMOTE とは、Chawla らによって提案されたオーバーサ ンプリング法の一つであり、k-最近傍法を基にしたアルゴ リズムである.任意の正例を指定し、類似度計算により最 近傍にある他の正例を特定する.そして、その二点を結ぶ 直線上に新たな正例を作成する.この作業を指定の回数繰 り返すことで正例の個数を増やし、データ間の偏りを減少 させる<sup>18)</sup>.式(14)は類似度計算の式である.

$$sim(m_a, m_i) = \sqrt{\sum_{j=1}^{n} (v_{a,j} - v_{i,j})^2}$$
(14)

ここで, *m* は任意の点, v は点 *m* における説明変数, n は 説明変数の数を表す.

本研究では、正例と負例の比が 1:4 なので、データが均 衡になるように、正例を4倍にオーバーサンプリングして 学習データを作成、分析を行った.図-6に不均衡データ整 形後の SVM および RF の地すべり地形評価図と的中率を示 す.データ整形前と比べて、全体の的中率は下がったが、 地すべり発生地点の抽出の精度は向上したため、オーバー サンプリングの効果と考えられる.



# 5. 熊本地震による土砂崩壊への適用

前章までの検討では、地すべりを引き起こした誘因を考 慮していない.そこで、熊本地震によって発生した土砂崩 壊地の地形的特徴を前章と同様に SVM, RF の二手法の機 械学習を用いて評価した.本章では前章までの分析で用い た説明変数のほかに、土砂崩壊の誘因として地震動強さを 考慮するため、QuiQuake<sup>19)</sup>の最大速度(PGV)分布を追加 した.また、土砂崩壊発生地点のデータには国土地理院の 土砂崩壊地分布図<sup>20)</sup>を利用した.これは、国土地理院が熊 本地震後に撮影した航空写真から、土砂崩壊地の分布を判 読したものである.対象は熊本県内 109741 メッシュである.

仙台市の分析と同様に,全データの10%を学習データと して分類器を構築し,残りの90%を予測データとして精度 を評価した.仙台市のデータと比べより負例の数が多い不 均衡データであったため,分析結果が負例に引っ張られて しまいどちらの手法でも土砂崩壊地の抽出の的中率が0% となった.そこで,この問題を解決するために SMOTE に よるオーバーサンプリングを行い,再度分析を実行した. 正例と負例の比率から,正例が160 倍となるようにデータ をオーバーサンプリングした.

データ整形後の機械学習の分析結果から作成した土砂崩

壊評価図および各手法の的中率を図-7に示す. 仙台市の分 析結果とは異なり, 各手法の的中率の差がみられ, SVM で は土砂崩壊地抽出の精度が, RF では全体の的中率が高かっ た.



SVM(SMOTE 適用後) 全体:93.6%, 土砂崩壞発生地点:60.0%



RF(SMOTE 適用後) 全体:98.6%, 土砂崩壊発生地点:25.4%

## 図-7 熊本地震後の土砂崩壊評価図

## 6. まとめ

本研究では、共分散構造分析とサポートベクターマシン、 ランダムフォレストの二種類の機械学習手法に基づき、宮 城県仙台市との地すべりの危険箇所の特徴を評価した.ま た、誘因を加えた分析として、熊本県内において熊本地震 により発生した土砂崩壊地の特徴を評価した.

機械学習を実行した結果,負例が正例に比べて極端に多い不均衡データであったため,分析による危険個所の見落としが多くなる問題が起こった.そこで,SMOTE によるオーバーサンプリングを行い,データを整形したところ,機械学習の精度を向上させることができた.機械学習の二手法を比較すると,仙台市の分析においては SVM と RF の全体精度,地すべり地形抽出精度ともにほぼ同程度であった.熊本県の分析では,土砂崩壊の誘因として PGV を考慮したことが影響し,機械学習の手法間で的中精度の違いが見られ,SVM の精度が RF よりも高かった.

# 参考文献

- 1) 新潟県土木部砂防課:新潟県中越地震と土砂災害, http://www.pref.niigata.lg.jp/HTML\_Article/864/780/chuet sujisin.pdf
- 国土交通省:災害・防災状況, http://www.mlit.go.jp/saigai/ saigai\_110311.html
- 国土交通省砂防部:平成26年8月豪雨による広島県で 発生した土砂災害への対応状況, http://www.mlit.go.jp/river /sabo/H26\_hiroshima/141031\_hiroshimadosekiryu.pdf
- 4) 豊田秀樹:共分散構造分析[R編],東京図書,2014.
- 5) 栗田多喜夫:サポートベクターマシン入門, http:// home.hiroshimau.ac.jp/tkurita/lecture/svm.pdf
- 6) 波部斉:ランダムフォレスト, https://www.slideshare. net/HitoshiHabe/ss-58784421
- 国土交通省:国土数値情報ダウンロードサービス, http://nlftp.mlit.go.jp/ksj/
- 若松加寿江,松岡昌志,久保純子,長谷川浩一,杉浦 正美:日本全国地形・地盤分類メッシュマップの構築, 土木学会論文集,No.759/1-67,213-232,2004.
- 防災科学技術研究所:地すべり地形分布図データベース,http://lsweb1.ess.bosai.go.jp/(2015年5月28日閲覧)
- 10) 小島尚人,大林成行,青木太:共分散構造分析を導入 した斜面崩壊危険箇所評価アルゴリズムの構築,土木 学会論文集,No.714/VI-56, pp.79-93,2002.
- 大林成行,小島尚人, Chang-Jo F.Chung:斜面安定性評価モデルの精度比較とその実用化への提案,土木学会論文集,No.630/VI-44, pp.77-89,1999.
- 12) 竹内一郎, 鳥山昌幸:サポートベクトルマシン(機械学 習プロフェッショナルシリーズ),講談社, 2015
- 13) 荒川正幹, 宮尾知幸, 船津公人:ドラッグライクネス モデルの構築とその可視化, Journal of Computer Aided Chemistry, Vol. 9, pp. 70-80, 2008.
- 波部斉: ランダムフォレスト,コンピュータビジョン とイメージメディア, pp. 1-8, 2012.
- 山岡 啓介: ランダムフォレスト, 映像情報メディ ア学会誌, Vol. 66, No. 7, pp. 573-575, 2012.
- 16) 小林寛武, 戸田航史, 亀井靖高, 門田暁人, 峯恒憲, 鵜 林尚靖:11 種類の fault 密度予測モデルの実証的評価, 電子情報通信学会論文誌 D, Vol.96, No. 8, pp. 1892-1902, 2013.
- Rpubs:SMOTE で不均衡データの分類, https://rpubs.com/ hoxo\_m/54954 (2017 年 7 月 4 日閲覧)
- 18) 亀井靖高,門田暁人,松本健一:Fault-proneness モデ ルへのオーバーサンプリング法の適用, http://www.empirical. jp/paperdb/papers/archive/102/102.pdf (2017年7月4日 閲覧)
- 19) 産業技術総合研究所: QuiQuake-地振動マップ即時推 定システム-, https://gbank.gsj.jp/QuiQuake/ (2016 年 10 月 11 日閲覧)
- 20) 国土地理院:土砂崩壞地分布, http://www.gsi.go.jp/BOUSAI/ H27-kumamoto-earthquake-index.html#cc