

卒業論文

地価形成に関わるオープンデータの ヘドニック・アプローチによる便益の予測

Visualizing Crime Factors and Improving the Accuracy of
Predictive Models Dealing with Imbalanced Data

富山県立大学 工学部 情報システム工学科

2120031 中島健希

指導教員 奥原 浩之 教授

提出年月: 令和5年（2023年）2月

目次

図一覧	ii
表一覧	iii
記号一覧	iv
第1章 はじめに	1
§ 1.1 本研究の背景	1
§ 1.2 本研究の目的	2
§ 1.3 本論文の概要	3
第2章 多様な要因を考慮したデータセットの作成	4
§ 2.1 サイバー空間からのデータ取得	4
§ 2.2 データセットに対する前処理	6
§ 2.3 説明変数の選定	8
第3章 ヘドニック・アプローチによる土地価格決定要因の分析	11
§ 3.1 ヘドニック・アプローチの理論モデル	11
§ 3.2 スパース推定	13
§ 3.3 Folium を用いた Web-GIS の開発	15
第4章 提案手法	18
§ 4.1 多様な要因を考慮したデータセットの作成	18
§ 4.2 不均衡なデータに対処した予測モデルの構築	21
§ 4.3 犯罪発生要因の可視化	24
第5章 数値実験並びに考察	27
§ 5.1 数値実験の概要	27
§ 5.2 実験結果と考察	28
第6章 おわりに	33
謝辞	34
参考文献	35

図一覧

2.1	Mapbox Studio	5
2.2	NAVITIME	5
2.3	犯罪発生地点を可視化した例	10
2.4	GIS によるホットスポット検出 [23]	10
3.1	データセットを作成するまでの流れ	13
3.2	Folium による Web-GIS 実装例	17
4.1	データセットを作成するまでの流れ	19
4.2	施設データに基づく説明変数	20
4.3	地図画像に基づく説明変数	20
4.4	予測モデルを作成する流れ	22
4.5	PR-AUC のグラフ例 [37]	23
4.6	混同行列と評価指数	23
4.7	要因マップを作成する流れ	25
4.8	犯罪発生要因マップの例	26
4.9	各要素の SHAP 値の例	26
5.1	予測モデルを検証する流れ	28
5.2	2020 年 9 月 1 日（左）と 2 日（右）の予測結果	28
5.3	予測モデルを可視化した結果	30
5.4	特定のグリッドセルの要因を可視化した結果	30

表一覧

3.1	代表的な GIS ソフトウェア	15
4.1	SHAP 値テーブルの例	25
5.1	データセットに含まれる罪種	28
5.2	検証用データによる予測結果	29
5.3	使用, および選択された説明変数一覧	32

記号一覧

以下に本論文において用いられる用語と記号の対応表を示す.

用語	記号
識別境界に直行している射影軸	\boldsymbol{w}
クラス間変動行列	\boldsymbol{S}_B
クラス内変動行列	\boldsymbol{S}_W
データセット内のクラス	C_n
クラスの平均値	\boldsymbol{m}_n
データセット内の多数クラス	C^{maj}
多数クラスのサンプル数	N^{maj}
データセット内の少数クラス	C^{min}
少数クラスのサンプル数	N^{min}
データセットの不均衡度	r
説明変数の集合	\boldsymbol{X}
学習済みのモデル	$\hat{f}(\boldsymbol{X})$
インスタンス i の説明変数の集合	\boldsymbol{x}_i
インスタンス i の予測値	$\hat{f}(\boldsymbol{x}_i)$
モデル $\hat{f}(\boldsymbol{X})$ の予測の期待値	$\mathbb{E}[\hat{f}(\boldsymbol{X})]$
インスタンス i の説明変数 $x_{i,j}$ の限界貢献度	$\Delta_{i,j}$
インスタンス i の説明変数 $x_{i,j}$ の SHAP 値	$\phi_{i,j}$
2 点 $(x_1, y_1), (x_2, y_2)$ の距離	D
2 点 $(x_1, y_1), (x_2, y_2)$ の緯度の差	D_y
2 点 $(x_1, y_1), (x_2, y_2)$ の経度の差	D_x
子午線曲率半径	M
卯酉線曲率半径	N
離心率	E
長半径	R_x
短半径	R_y
道路ネットワークにおける平均次数	k

はじめに

§ 1.1 本研究の背景

不動産価格は「一物四価」とも言われ、実勢価格、公示地価、路線価、固定資産税評価額の4種類が存在する。実勢価格は市場において現実に成立した価格であり、一般的な売買の価格を決定する際に参考とされるが、売主と買主の個別の事情によって価格が決定されるため、必ずしも適正な価格が形成されているとは限らない。次に、公示地価は、地価公示法に基づいて評価された毎年1月1日時点の1m²当たりの価格であり、毎年3月頃に公示され、一般の土地の取引価格の指標とされる。また、路線価は土地に対する相続税や贈与税算定のための基礎価格として、国税庁が発表する価格のことであり、公示地価の約80%を目安に決定され、毎年1月1日を評価時点として、地価公示価格・基準地価・売買事例などをもとに算出される。最後に、固定資産税評価額は、各市区町村が算定する固定資産税の基準となる評価額であり、主に固定資産税計算時に用いられ、公示地価の約70%が目安となっている。

不動産の価値評価に関して、国土交通省の不動産鑑定評価基準では次のように記されている。

不動産の価格は、多数の要因の相互作用の結果として形成されるものであるが、要因それ自体も常に変動する傾向を持っている。したがって、不動産の鑑定評価を行うに当たっては、価格形成要因を市場参加者の観点から明確に把握し、かつ、その推移及び動向並びに諸要因間の相互関係を十分に分析して、前記三者（不動産の効用、相対的稀少性、不動産に対する有効需要）に及ぼすその影響を判定することが必要である。〔国土交通省（2002）p.6〕

具体的に価格形成要因は、一般的要因、地域要因及び個別的要因に分けられるとされている。

多くの要因ゆえに、一般的に住宅地地価を含む不動産はユニークな商品であるという側面を持つ。すなわち、自然的特性としての不動産、再生産不可能性、土地やその形状といった個別性、加えて人的特性としての用途の多様性ゆえに、同一の不動産というものは存在せず、それぞれは完全に差別化された財といえるのである。要因それぞれの重要度を再整理したうえで定量化するために、本論ではヘドニック・アプローチ（hedonic approach）を採用する。ヘドニック・アプローチでは住宅地の価値を地域のアメニティ、環境質、利便性、規模等のような「特性（characteristics）」の合成物であると考え、この場合、共通の客観的性質を示す特性のレベルに依存して住宅地地価が決定されることが考えられることから、共通の客観的性質に基づく金額換算が可能となる。概念的には、観測・非観測な情報分析を通じて、住宅地需要者が一定の予算制約のもとで効用最大化が図れる土地を選択し、一方住宅地供

給者は自らの利潤最大化が図れる場合に住宅地の供給を行うというものである。そうした結果として成立する均衡価格を想定するのがヘドニック・アプローチである。このようなアプローチは不動産市場分析でよく使われる手法である一方、社会資本等の非市場の経済評価や、日本銀行調査統計局（2007）のように品質調整法として活用されることも多い。

§ 1.2 本研究の目的

ヘドニックアプローチは、物価指数の作成における品質調整方法の一つであり、商品の品質をその商品の持つ個々の機能・性能に分解可能であるとの仮定に基づき、各機能・性能と価格の関係を回帰分析を通じて推計する手法である。この手法は、不動産のように多様な特性が価格形成に関与する場合に有効であり、住宅地の価値が地域の特性や環境の質に依存するという考え方に基づいて、各特性の価格への影響を評価することが可能となる。しかし、ヘドニックアプローチでは、説明変数間の多重共線性の問題や、変数の欠落によるバイアスの発生が課題とされている。

本研究では、これらの課題を克服するために、スパース推定の一つであるアダプティブエラスティックネット（Adaptive Elastic Net, AEN）を導入する。AENは、多数の説明変数の中から意味のある変数のみを選択し、それ以外の変数のパラメータを0とする特性を持つため、変数選択の自動化を可能にする。さらに、AENは「グループ効果」と呼ばれる多重共線性に対する頑健性や、変数選択の適正性と係数の漸近的な不偏性を同時に担保する「オラクル性」を有している。これにより、変数間の多重共線性や、変数の欠落によるバイアスの問題を克服できる可能性がある。

本研究を行う理由は、不動産価格の決定要因を明確化することが、適正な不動産評価や価格予測の精度向上に寄与すると考えたからである。特に、土地取引の価格は、地域経済や投資判断に大きな影響を与えるため、より客観的で信頼性の高い評価手法の開発が求められている。また、不動産価格は地域社会の経済活動の指標の一つでもあり、政策立案や不動産市場の安定化にも重要な示唆を与える。ヘドニックアプローチを用いた不動産価格の分析は、不動産取引の透明性を高め、より公正な取引環境を構築するためにも必要不可欠な手法であるといえる。

本研究の重要性は、従来の最小二乗法を用いたヘドニック分析の限界を克服する新たな手法を提示する点にある。特に、AENを用いたスパース推定による変数選択の自動化は、多数の説明変数が存在する場合でも、重要な変数のみを適切に選択できる可能性があるため、予測精度の向上とともに、モデルの解釈性を高める効果が期待される。本研究は、不動産の価格形成メカニズムの理解を深めるとともに、価格評価の公正性と透明性を高める新たなアプローチを提供するものであり、不動産市場の発展に寄与することが期待される。

§ 1.3 本論文の概要

本論文は次のように構成される。

- 第1章** 本研究の背景と目的について説明した。背景では、特に欧米における地理的犯罪予測の歴史と事例について述べた。目的では、わが国における地理的犯罪予測の課題について述べ、本研究の意義について述べた。
- 第2章** 地理的犯罪予測の概要と、その手法についてそれぞれ述べる。また、犯罪が発生するリスクについて述べる。さらに、地理的犯罪予測には欠かせないGISについて、その概要を述べる。
- 第3章** 不均衡なデータに対するアプローチと、機械学習によって作成されたモデルを解釈する手法について述べる。また、さまざまな要因を考慮するため、サイバー空間から多様なデータを取得し、処理する方法について述べる。
- 第4章** データセットを作成し、不均衡に対処して犯罪発生予測モデルを作成する。さらに、その予測モデルに解釈手法を適用し、犯罪が発生する要因を可視化するまでの流れを説明する。
- 第5章** 実際の犯罪発生データを用いて、第4章で述べた手法で、犯罪発生予測モデルを作成し、その予測精度を検証する。また、解釈手法によって可視化された要因が妥当なものであるかを確認する。
- 第6章** 本論文における前章までの内容をまとめつつ、本研究で実現できたことと今後の展望について述べる。

多様な要因を考慮したデータセットの作成

§ 2.1 サイバー空間からのデータ取得

土地価格の変動には、数多くの要因が考えられる。そのため、土地価格を予測するモデルを作成するためには、それらを表現する説明変数を多く考慮する必要がある。しかし、我々が一般に取得できるデータ、すなわちオープンデータには、そのアクセスに限界がある。実際に、日本で公開されているオープンデータの数、世界で最も公開されている台湾と比較して、約 67.7 % である [28]。国勢調査の結果など、統計的なデータは比較的公開されているものの、土地価格の要因として重要視される地理的なデータ、たとえば、特定の施設の位置などといったものは、依然として取得が容易ではない。そこで、本研究では、地理的なデータを地図画像やナビゲーションサービスから取得し、補うこととした。

地図画像の取得

地図画像は、その場所やその周囲の地理的な特徴を表す重要なデータである。そこで、本研究では、Mapbox から取得した地図画像から説明変数を抽出している。Mapbox は、機能 14 やデザインを自由にカスタムして、地図を自身の Web ページやアプリに埋め込むことができるサービスである。さまざまな API を公開しており、住所などから緯度・経度を算出する Geocoding API、ルートを検索する Directions API などがあるが、本研究では、地図をベクター画像として取得できる Mapbox Static Tiles API を用いて、地理的なデータを取得する。

Step 1: Mapbox Studio 上で、カスタムマップを作成する

Mapbox Studio では、地図上にあるさまざまな要素の色や表示の有無を自由に変更することができる。

Step 2: 緯度と経度から、取得するタイルを算出する

Mapbox Static Tiles API では、地球上のすべての範囲を正方形で仕切ったタイルごとに地図画像を取得できる。すなわち、緯度と経度から、特定のタイルを一意に決定することができる。対象の緯度と経度 (lat, lng) が含まれるタイル (X, Y) は、次のように算出できる。

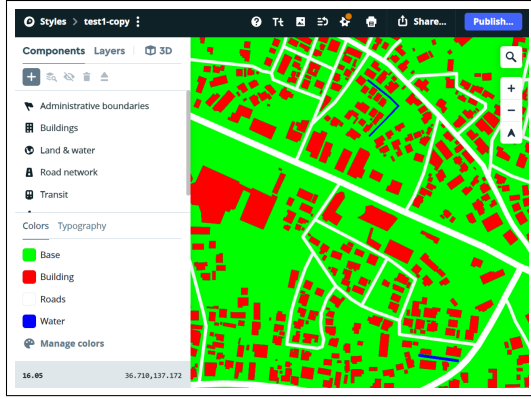


図 2.1: Mapbox Studio



図 2.2: NAVITIME

$$X = \lfloor \frac{lon + 180}{360} * 2^z \rfloor \quad (2.1)$$

$$Y = \lfloor \frac{\log_e \tan \left(lat \frac{\pi}{180} + \frac{1}{\cos \left(lat \frac{\pi}{180} \right)} \right)}{\pi} 2^{z-1} \rfloor \quad (2.2)$$

ここで、 $\lfloor x \rfloor$ は、 $n \leq x < n+1$ を満たす整数 n を表す。また、 z はズームレベルである。たとえば、 $z = 17$ では1ピクセルあたり 1.194m、 $z = 18$ では1ピクセルあたり 0.597m の地図画像を取得できる。Mapbox Static Tiles APIで取得できる地図画像の大きさは 512×512 であるため、 $z = 17$ では一辺が約 611m、 $z = 18$ では約 306m である。

Step 3: Mapbox Static Tiles APIを用いて、地図画像を取得する

以上により、タイル X, Y 、およびズームレベル z を算出・決定したら、Mapbox Static Tiles APIとして指定されている URL に、それらをパラメータとして GET リクエストを行う。レスポンスされたデータはバイト列であるため、1つの座標に RGB 値を格納する 3次元配列に変換を行えば、画像として処理することができる。

施設データの取得

特定の施設やその近くは、犯罪の発生の要因となる可能性がある。施設データを取得できるサービスとして、Google Maps APIが存在するが、無料で取得できる数に制限があるほか、たとえば遊園地や水族館など、レジャー施設としてジャンル分けできるものに対し

て、「レジャー施設」と検索しても、それらを網羅できるとは限らない点で、採用しなかった。そこで、ナビゲーションサービスのひとつである「NAVITIME」から施設データをスクレイピングして取得することとした。

スクレイピングとは、データを収集した上で利用しやすいように加工をすることである。特に、Web 上から必要なデータを取得することを、Web スクレイピングと呼ばれている。スクレイピングと似ている意味の言葉にクロールリングがあり、スクレイピングとは違い、これは、単に Web 上のデータを収集することを意味する。データを活用するために、使いやすく抽出や加工をしたりするのがスクレイピングの特徴である。

BeautifulSoup4 とは、Web サイト上の HTML から、必要なデータを抽出することができるライブラリである。Beautifulsoup4 でスクレイピングする際、最初に対象の Web ページから HTML を取得する必要がある。HTML を取得する方法として、同じく Python のライブラリである、Requests の get 関数などがある。上記の方法によって取得された HTML テキストを、BeautifulSoup4 の BeautifulSoup 関数に渡すことで BeautifulSoup オブジェクトを作成ができ、そのオブジェクトから要素検索をすることで必要な情報を抽出する。

NAVITIME は、施設のジャンルごと、さらには都道府県ごとに一覧となつて表示される。

§ 2.2 データセットに対する前処理

土地価格の変動には、数多くの要因が考えられる。そのため、土地価格を予測するモデルを作成するためには、それらを表現する説明変数を多く考慮する必要がある。しかし、生のデータには欠損値や外れ値が含まれていることが多く、そのままではモデルの精度が低下する恐れがある。本節では、回帰分析におけるデータセットの前処理手法について説明する。

データの前処理の概要

前処理は、モデルの精度向上や学習の安定性を確保するために不可欠な手順である。データセットの前処理では、以下の手順を実施する。

Step 1: 欠損値の処理 データセットには、しばしば欠損値が含まれている。欠損値の処理には、以下のような方法がある。

- **削除**：欠損のあるデータポイントを削除する。ただし、サンプルサイズが減少するリスクがある。
- **補完**：欠損値を他の値で補完する方法。平均値、中央値、最頻値での補完や、k-近傍法 (kNN) を用いた補完がある。
- **モデルを用いた補完**：機械学習モデルを用いて、欠損部分を予測する方法もある。

Step 2: 外れ値の処理

外れ値は、モデルの予測精度を低下させる要因となる。外れ値の検出方法には、以下の手法がある。

- **四分位範囲（IQR）による検出**：四分位範囲（IQR）の 1.5 倍を超えるデータポイントを外れ値とみなす方法.
- **標準偏差を用いた検出**：平均からの標準偏差が一定の範囲を超えるデータを外れ値とする方法.
- **視覚的検出**：散布図や箱ひげ図を用いて視覚的に外れ値を確認する方法.

外れ値は、そのまま維持するか、削除するか、または他の値に変換する（Winsorization）かのいずれかを選択する.

Step 3: カテゴリ変数のダミー変数化

ダミー変数とはカテゴリカルデータを「0」または「1」の数値データに変換した変数のことである. カテゴリ変数は、機械学習モデルでは直接使用できないため、数値データに変換する必要がある. 代表的な方法は、**ワンホットエンコーディング**である.

- **ワンホットエンコーディング**：カテゴリ変数を 0 または 1 の二値変数に変換する方法. たとえば、{A, B, C} という 3 つのカテゴリがある場合、これを [1, 0, 0], [0, 1, 0], [0, 0, 1] といったベクトルに変換する.
- **ラベルエンコーディング**：カテゴリ変数に対して整数を割り当てる方法. ただし、この手法はカテゴリの大小関係が生じてしまうため、回帰分析には不向きな場合がある.
- **3 つ以上のカテゴリを持つ変数のダミー変数化**：3 つ以上のカテゴリを持つ変数の場合、ワンホットエンコーディングにより、各カテゴリに対応するダミー変数を作成する. 例えば、{A, B, C, D} というカテゴリ変数がある場合、これを {[1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0], [0, 0, 0, 1]} のように変換することができる.
- **ダミー変数の落とし込み（ダミー変数の落とし込み問題）**：ダミー変数化の際、冗長な変数が生じることを避けるため、1 つのカテゴリ変数を削除することが一般的である. 例えば、4 つのカテゴリ変数がある場合、3 つのダミー変数を作成し、残り 1 つを基準カテゴリとして扱う. この基準カテゴリは、回帰分析において「参照カテゴリ」となり、他のカテゴリとの相対的な影響を示す.
- **注意点**：ダミー変数を多く生成しすぎると、モデルの計算負担が大きくなることがあるため、カテゴリ数が多い場合には、次元削減技術（例えば、主成分分析（PCA）など）を検討することが望ましい.

Step 4: 標準化

1. **Z スコアによる標準化** 各変数の値から平均を引き、標準偏差で割ることで、平均が 0、標準偏差が 1 となるように変換する. これにより、異なるスケールの変数を均一な基準に揃えることができる. 標準化の数式は次の通りである.

$$z = \frac{x - \mu}{\sigma} \quad (2.3)$$

ここで, x は元の値, μ は平均, σ は標準偏差である.

2. **Min-Max スケーリング** 変数の最小値を 0, 最大値を 1 に変換する方法である. すべてのデータが 0 から 1 の範囲に収まるため, ニューラルネットワークのようなモデルでよく用いられる. スケーリングの数式は以下の通りである.

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad (2.4)$$

ここで, x は元の値, x_{\min} は最小値, x_{\max} は最大値である.

Step 5: データの分割

データセットを, モデルの学習用 (訓練データ) と性能評価用 (テストデータ) に分割する. 分割比率は, 80:20 や 70:30 が一般的である.

前処理の重要性

前処理は, モデルの予測精度や安定性に大きな影響を与えるため, 回帰分析の成功にとって極めて重要である. 欠損値や外れ値がある状態でモデルを学習させると, モデルのパフォーマンスが低下することがあるため, 適切な前処理が求められる. また, カテゴリ変数をダミー変数に変換しないままモデルに投入すると, 予測誤差の増大を招く可能性がある. そのため, データセットの内容を把握し, 必要な前処理を慎重に行うことが不可欠である.

§ 2.3 説明変数の選定

2.1 節で述べた RTM のように, 地理的犯罪予測において, GIS は欠かせない重要な要素となっている. 本節では, その GIS について述べるとともに, GIS の地理的犯罪予測への応用について説明する.

GIS

GIS とは, 位置に関する情報を持ったデータを総合的に管理・加工し, 地理的な位置とデータを結び付けることによって, 視覚的に表示することである. これにより, 空間データの高度な分析や迅速な判断を可能にする技術である. 地理空間データの例として, 主題図 (土地利用図, 地質図, ハザードマップ等), 都市計画図, 地形図, 地名情報, 台帳情報, 統計情報, 空中写真, 衛星画像などが挙げられる. GIS はいくつかの優れた特徴を持つため, さまざまな用途で用いられている. その中でも代表的なものとして, 次のような 4 つの特徴と 3 つの利点が挙げられる [22].

GIS における 4 つの特徴

- データの可視化

2D や 3D、アニメーションなど多様な表現方法で地図上にデータを可視化することによって、数値のみを見るだけでは気づくことができないようなデータの傾向やデータ間の関連性など様々な情報を一目で把握できるようになる。

- データ間の関係性の把握

複数のデータによる重ね合わせを行うことでデータ間の対比が容易になり、データ同士の関係性を直感的に把握することができる。また、地図上の位置関係からデータを特定することによって、定量的な情報を把握することができる。

- データの統合と分析

位置情報をキーとしてそれぞれ別の特徴を持つデータを統合したり、複数のデータを重ね合わせて分析することによって、独自のアプローチで課題の解を導き出すことができる。

- データの作成と更新

今日の日本では新しいビルの建設や合併による行政界の変更など、現実世界が日々変化しているため地理情報データもそれらに基づいて定期的に更新していく必要がある。この際、GIS を利用することでデータの作成・更新における負担が軽減され、常に鮮度を保ってデータの管理・提供を行うことができる。

GIS を用いることの 3 つの利点

- 業務効率化によるコスト削減

GIS は日常の業務を最適化するために幅広く利用されている。紙地図から GIS を利用したデジタルな地図へ移行することで、現地調査や設備管理、統計分析などを行う際の物質的な制約が減少し、より簡単かつ効率的な作業の実現が可能になる。これにより作業時間や人員など業務にかかるコストを大幅に削減することができるという利点を持つ。

- 最適な意思決定

組織における活動において、場所に関する意思決定を正しく行うことは成功のための重要な要素の 1 つといえる。場所に関する決定に視覚を用いることができるという利点を持つ GIS は店舗の出店場所や配送ルート、避難地域・経路、天然資源の採取地点など多様な分野で最適な場所を策定するために使用されている。

- コミュニケーションの向上

GIS を使用することでさまざまな表現方法を用いて位置情報を地図上に可視化することができる。可視化された位置情報は状況を効果的に伝え、的確な理解を促すことができ、グループや組織間、社会におけるコミュニケーションの向上を図ることが可能になるという利点を持つ。

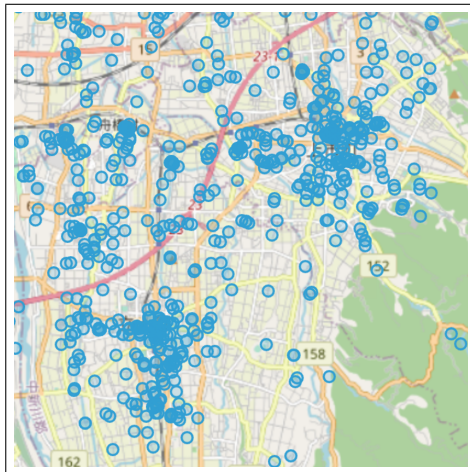


図 2.3: 犯罪発生地点を可視化した例

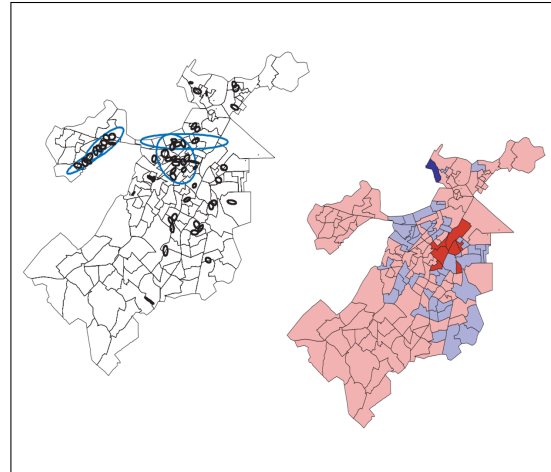


図 2.4: GIS によるホットスポット検出 [23]

地理的犯罪予測への GIS の応用

犯罪にはホットスポットと呼ばれる概念が存在する。ホットスポットは、犯罪が時空間的に集積して発生する現象であり、すなわち、ホットスポットは何らかの脆弱性が存在し、その同定は、将来に発生する犯罪を予測するためにも有用である。このようなホットスポットを分析するうえで、犯罪発生状況を地図上に可視化するクライムマッピングは重要となる。

特に、ホットスポットを特定するためによく用いられるカーネル密度推定は、GIS の発達により容易になったといっても過言ではない。このように、GIS が発達し、誰でも容易に用いることができるようになったことで、ホットスポットの同定が行えるようになり、ホットスポットに関する理論や分析が、地理的犯罪予測の基礎となったといえる [7]。また、地理的犯罪予測の発展には、警察機関の動向も大きく関わっている。

1990 年代に米国の警察が、GIS を日常的に使用するようになったことで、地理的犯罪予測に関する研究が一気に加速したと述べている [24]。その例として、当時のニューヨーク市警察が取り入れた「コムスタット」が挙げられる。これは、犯罪データを週単位で集計して地図上で可視化し、どのような戦略を行うのか、その立案を行うものである。この GIS を取り入れた警察活動は、地理的犯罪予測と、その結果に基づく警察活動の基盤となった。

ヘドニック・アプローチによる土地価格決定要因の分析

§ 3.1 ヘドニック・アプローチの理論モデル

本研究では、Epple（1987）の議論に基づき、住宅地地価の諸特性を取引する暗黙的な市場を想定したうえで、市場均衡価格曲線としてのヘドニック関数を導出する。以下では、需要サイドおよび供給サイドにおける行動を定式化し、市場均衡に至る過程を示す。

需要サイドの行動

地価を形成する特性の n 次元ベクトルを $\mathbf{z} = (z_1, z_2, \dots, z_n)$ 、ヘドニック価格関数を $p(\mathbf{z}) = p(z_1, z_2, \dots, z_n)$ とする。また、住宅地需要者の効用関数を $U(\mathbf{z}, x; \boldsymbol{\alpha})$ と定義する。ここで、 x はニューメレール（価値尺度財）、 $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)$ は需要者個人のテイストパラメータのベクトルである。需要者の所得を y とした場合、予算制約式は以下のように表される。

$$y = p(\mathbf{z}) + x \quad (3.1)$$

同時分布関数を $F(y, \boldsymbol{\alpha})$ と表す。この予算制約式のもと、需要者が \mathbf{z} および x について効用最大化行動を取ると、次式に定式化される。

$$\max_{\mathbf{z}, x} U(\mathbf{z}, x; \boldsymbol{\alpha}) \quad (3.2)$$

$$\text{s.t. } y = p(\mathbf{z}) + x \quad (3.3)$$

この場合、最適化のための1階条件（FOC）は以下の式で表される。

$$p_{\mathbf{z}} = \frac{U_{\mathbf{z}}(\mathbf{z}, y - p(\mathbf{z}); \boldsymbol{\alpha})}{U_x(\mathbf{z}, y - p(\mathbf{z}); \boldsymbol{\alpha})} = h(\mathbf{z}, y - p(\mathbf{z}); \boldsymbol{\alpha}) \quad (3.4)$$

ここで、 $p_{\mathbf{z}}$ はヘドニック価格関数の1階微分のベクトルであり、 $U_{\mathbf{z}}$ および U_x はそれぞれ特性ベクトル \mathbf{z} およびニューメレール x の1階微分を示す。

需要者の効用水準が u の下でのビッド関数（bid function）を $\theta(\mathbf{z}; u, y)$ とすると、 $U(\mathbf{z}, y - \theta) = u$ が成立し、これを微分することで次式が得られる。

$$\frac{\partial \theta}{\partial z_i} = \frac{U_{z_i}}{U_x} > 0 \quad (3.5)$$

$$\frac{\partial^2 \theta}{\partial z_i^2} = \frac{U_{z_i}^2 U_{xx} - 2U_{z_i} U_x U_{z_i x} + U_x^2 U_{z_i z_i}}{U_x^3} < 0 \quad (3.6)$$

すなわち., ビッド関数は増加する凹関数である. 需要者の効用は., ヘドニック関数とビッド関数の接点において最大化されるため., 次の式が成立する.

$$\theta(\mathbf{z}^*; u^*, y) = p(\mathbf{z}^*) \quad (3.7)$$

$$\frac{\partial \theta}{\partial \mathbf{z}}(\mathbf{z}^*; u^*, y) = p_{\mathbf{z}}(\mathbf{z}^*) \quad (3.8)$$

図的には., ヘドニック関数がビッド関数のエンベロップ・カーブとなる.

供給サイドの行動

次に., 供給者の行動を定式化する. 供給者は自らの供給行動を決定する際., 住宅地地価を所与として利潤 π を最大化するように特性の束 $\mathbf{z} = (z_1, z_2, \dots, z_n)$ を選択する. 利潤関数は以下のように表される.

$$\max_{\mathbf{z}, M} \quad \pi = p(\mathbf{z})M - C(M, \mathbf{z}; \boldsymbol{\beta}) \quad (3.9)$$

ここで., M は供給する住宅地の数., $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_n)$ は供給者を特徴づけるパラメータベクトルであり., その分布関数を $G(\boldsymbol{\beta})$ とする. また., $C(M, \mathbf{z}; \boldsymbol{\beta})$ は供給者の費用関数である. この場合., 利潤最大化の 1 階条件は以下ようになる.

$$p_{\mathbf{z}} = C_{\mathbf{z}}(M, \mathbf{z}; \boldsymbol{\beta}) \quad (3.10)$$

$$p(\mathbf{z}) = C_M(M, \mathbf{z}; \boldsymbol{\beta}) \quad (3.11)$$

ここで., 供給者は各特性の限界的な価値と土地 1 単位当たり特性の限界費用が等しくなるように供給活動を行う. また., 土地の市場価値は供給限界費用に等しくなる. 供給者のオファー関数 (offer function) を $\phi(\mathbf{z}, \pi)$ とすると., 以下が成立する.

$$\phi_{\mathbf{z}} = \frac{C_{\mathbf{z}}}{M} > 0 \quad (3.12)$$

$$\phi_{\pi} = \frac{1}{M} > 0 \quad (3.13)$$

すなわち., オファー関数は増加する凸関数である. 市場均衡は以下の式を満たす.

$$p(\mathbf{z}^*) = \phi(\mathbf{z}^*, \pi^*) \quad (3.14)$$

$$p_{\mathbf{z}}(\mathbf{z}^*) = \phi_{\mathbf{z}}(\mathbf{z}^*, \pi^*) \quad (3.15)$$

このように., ヘドニック関数は需要者のビッド関数と供給者のオファー関数が市場均衡価格を挟んで接する形で決定される.

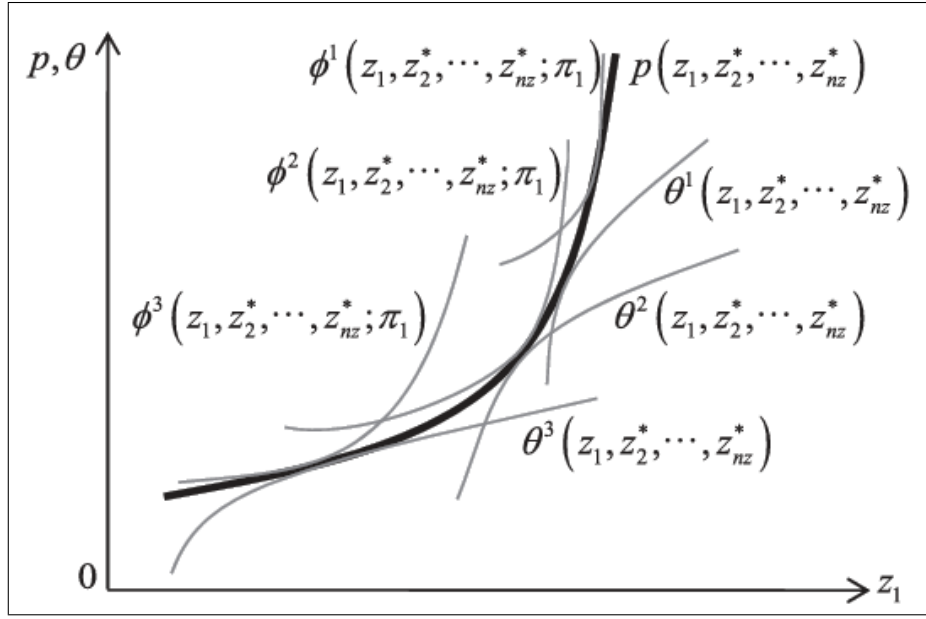


図 3.1: データセットを作成するまでの流れ

ヘドニック関数と市場均衡

ヘドニックアプローチでは、特性 $z = (z_1, z_2, \dots, z_n)$ を持つ住宅地に対する需要と供給が合致する点で市場均衡価格が決定される。この価格は需要者サイドの分布 $F(y, \alpha)$ と供給者サイドの分布 $G(\beta)$ に依存して決定される。

しかしながら、 $F(y, \alpha)$ や $G(\beta)$ が未知であるため、一般的には $p(z)$ も未知であり、需要サイドと供給サイドを同時推定することで市場均衡価格を導出する必要がある。この場合、同時方程式バイアスや関数型の問題が生じることが、清水・唐渡（2007）で指摘されている。

§ 3.2 スパース推定

スパース推定は、モデルの複雑さを抑制しつつ、予測性能を向上させるための手法である。特に高次元データにおいては、すべての説明変数をモデルに含めるのではなく、一部の重要な変数のみを選択することが望ましい。このような背景から、スパース推定は回帰分析において重要な役割を果たす。

スパース推定の代表的な手法

スパース推定の代表的な手法には、以下のものがある。

Lasso 回帰

Lasso (Least Absolute Shrinkage and Selection Operator) 回帰は、目的関数に ℓ_1 正則化項を加えることで、不要な変数の係数を 0 にする手法である。これにより、変数選択が自動的に行われ、モデルの解釈性が向上する。

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p |\beta_j| \quad (3.16)$$

ここで、 y_i は目的変数、 x_i は説明変数、 β は回帰係数、 λ は正則化パラメータである。 λ の値が大きいほど、より多くの変数の係数が 0 になる。

Ridge 回帰

Ridge 回帰は、目的関数に ℓ_2 正則化項を加える手法である。 ℓ_2 正則化は、すべての係数を小さな値に抑制するが、Lasso のように係数を 0 にはしない。

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \quad (3.17)$$

Ridge 回帰は、多重共線性の問題を緩和する効果があり、すべての変数をモデルに残しながら過学習を防ぐ役割を果たす。

Elastic Net

Elastic Net は、Lasso と Ridge の正則化項を組み合わせた手法である。 ℓ_1 と ℓ_2 の両方の効果を取り入れることで、変数選択と多重共線性の両方の問題に対処できる。

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2 \quad (3.18)$$

ここで、 λ_1 と λ_2 はそれぞれ ℓ_1 と ℓ_2 の正則化パラメータである。Elastic Net は、Lasso と Ridge の利点を同時に享受できる手法として広く用いられている。

スパース推定の効果と利点

スパース推定には、以下の利点がある。

- **変数選択の自動化**：Lasso 回帰や Elastic Net は、不要な変数を自動的に除去するため、モデルの解釈性が向上する。
- **過学習の防止**：Ridge 回帰は、すべての変数をモデルに残しながら係数を小さく抑えるため、過学習のリスクを低減できる。
- **高次元データへの対応**：Lasso や Elastic Net は、変数の数が観測データ数よりも多い場合でも使用可能である。

グループ効果

グループ効果 (Group Effect) は、関連する変数がグループとしてまとめて選択される現象を指す。Lasso では一部の変数のみが選択されるが、関連する変数を同時に選択したい場合には、**Group Lasso** が用いられる。Group Lasso は、変数を事前に定義されたグループに分割し、グループごとに ℓ_2 ノルムを適用しながら、全体として ℓ_1 ノルムの正則化を行う。

$$\underset{\beta}{\text{minimize}} \quad \frac{1}{2n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{g=1}^G \|\beta^{(g)}\|_2 \quad (3.19)$$

ここで、 $\beta^{(g)}$ はグループ g に属する回帰係数ベクトルである。これにより、関連する変数がグループ単位で選択される。

オラクル性

オラクル性 (Oracle Property) とは、スパース推定手法が真のモデル構造を正しく推定できる性質を指す。オラクル性を満たす手法は、真に有効な変数を正しく特定し、不要な変数の係数を 0 にできる。Lasso は特定の条件下でオラクル性を満たすとされるが、Elastic Net や Ridge 回帰は必ずしもオラクル性を満たすわけではない。これを達成するためには、非ゼロの係数が十分に大きいことや、サンプルサイズが十分に大きいことが必要とされる。

スパース推定の応用事例

スパース推定は、以下のような分野で活用されている。

- **金融工学**：リスク管理やポートフォリオ最適化において、不要な資産を選択的に除外するために使用される。
- **医療分野**：ゲノムデータ解析において、重要な遺伝子を特定するために Lasso や Elastic Net が活用されている。
- **マーケティング**：顧客の購買行動を予測するためのモデルにおいて、影響力の大きい要因を特定するために利用されている。
- **不動産評価**：土地価格の評価モデルにおいて、不要な説明変数を除外し、重要な要因を特定するために用いられる。

§ 3.3 Folium を用いた Web-GIS の開発

2.2 節, 2.3 節では、すべての GIS に対して一般的な名称として GIS と表記したが、GIS には動作するプラットフォームや形態、提供している団体によって複数種類のソフトウェアが存在する。代表的な GIS アプリケーションを表 3.1 に示す。また、本節では、Web アプリケーションとして World Wide Web 上で機能する GIS を Web-GIS と表記するものとする。

html 形式で記述され、Web-GIS は多様なプラットフォーム上で動作する GIS の中でも代表的なフォーマットであるといえる。計算機を用いて広く利用することができ、処理の大部分は html が置かれているサーバ上で行われることで参照する端末自体のスペックに依存しづらいことから、利用者も多い。また、html 形式で実装されるがゆえに作成者が直接作成する方法のほか、いくつかのプログラミング言語を用いて自動的に生成することも可能である。

本節では、以上のような特徴を有する Web-GIS の開発において、現在、一般に広く用いられているプログラミング言語の一つである Python のライブラリを用いた方法を解説する。また、Web-GIS において実装することができる代表的な機能とその役割を解説する。

Web-GIS の実装方法

表 3.1: 代表的な GIS ソフトウェア

GIS ソフトウェア	無料	オープン ソース	Windows	MacOS	Linux	BSD	Unix	Web
Microsoft MapPoint	×	×	○	×	×	×	×	○
ArcGIS	×	×	○	○	×	×	○	○
GRASS GIS	○	○	○	○	○	○	○	○
QGIS	○	○	○	○	○	○	○	○
MapInfo	×	×	○	×	×	×	×	○
TNTmips	×	×	○	○	○	×	○	×

Python を用いた Web-GIS の開発には「Folium」という Python 用のライブラリを用いる。Folium を用いてメソッドに対して初期位置，ベーススタイル，初期縮尺などを引数として与えてプログラムを実行することで Web-GIS のベースとなるマップを表示する html が自動的に生成される。ベーススタイルとして指定することの出来るマップスタイルには代表的な例として以下のようなものがある [?].

- CartoDB (positron and dark_matter)
- OpenStreetMap
- Mapbox Bright
- Cloudmade
- Mapbox

Folium による Web-GIS の開発はこのベースマップに対して Folium のライブラリ内に含まれる様々なメソッドを用いることで Web-GIS における各種機能や実際に表示した情報を追加するという形で行われる。ここからは，Folium によって実装することの出来る各種機能とその内容について代表的なものを取り上げる。

ベースマップの切り替えとレイヤの重ね合わせ

2.3 で言及したような GIS によるデータの重ね合わせは FeatureGroup 関数によるレイヤの作成とそれらを制御する LayerControl メソッドによって実現される。また，レイヤとは Web-GIS 上でマップの重ね合わせを行う際のそれぞれの層のことを表す。

まず，FeatureGroup 関数によってベースマップとは別のマップ（レイヤ）を任意の個数作成する。次に，後述する様々なメソッドを用いてベースマップや各レイヤに対して各種機能や情報を追加する。

この時，常に表示され，レイヤの切り替えに左右されないようにする必要のある情報に関してはベースマップに，それ以外の情報に関しては切り替えによって表示したい各レイヤに追加するようにする。最後に，LayerControl メソッドを用いてレイヤを管理する機能を Web-GIS に付与することによって，html を生成した際に自由にレイヤを切り替える機能を持ったものが生成される。

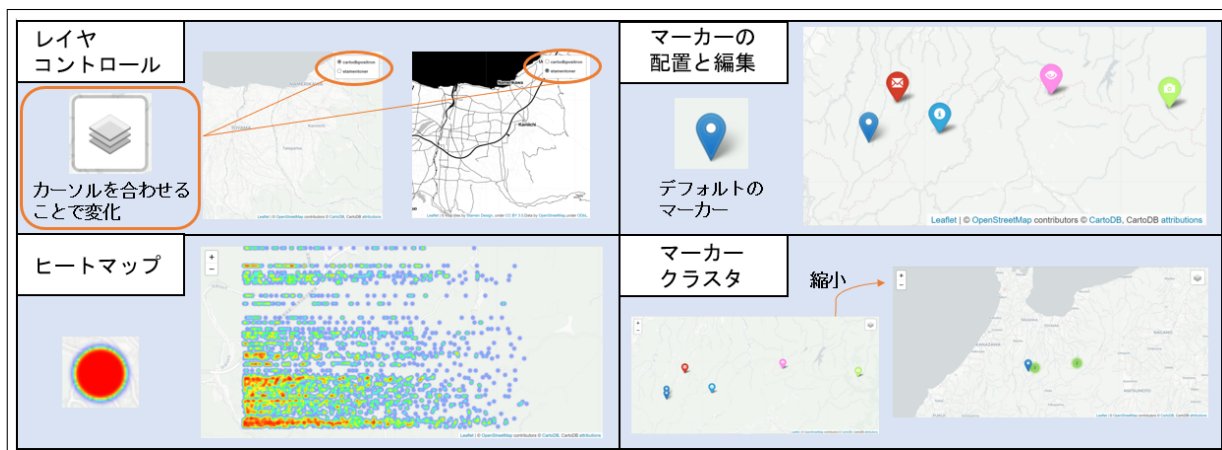


図 3.2: Folium による Web-GIS 実装例

マーカーを置く

`folium.Marker` メソッドに対して引数としてマーカーを置く位置の座標を与えることで地図上の任意の位置にマーカーを立てることができる。なお、マーカーは Folium 内に組み込まれているものの中から色やアイコンのマークを自由に切り替えて使用することができるほか、`CustomIcon` 関数によってアイコンのマークを自作し、独自のマーカーとして使用することもできる。また、任意のマーカーに対して `popup` 機能を追加し、テキストを付与しておくことでマーカーをクリックした際にポップアップとしてテキストが表示されるようになる。

ヒートマップを描く

ヒートマップとは、二次元データの数値を色やその濃淡で表したものである。広義におけるヒートマップは「マップ」と付いてはいるが必ずしも地図で表現する訳ではなく、テーブルを値で色分けしたものなど数値データを色分けによって可視化したものすべてがこれにあたる。ただ、Web-GIS におけるヒートマップは地図上にプロットされた色の濃淡で数値の大小を示すものである。

また、地図の色分けによって数値の大小を表現する方法として、ヒートマップとは様式が異なるものとして、コロプレス図(階級区分図)がある。コロプレス図とは、例えば地図を都道府県ごとに境界線で分けて、各都道府県における統計データの大きさによって色分けするなどのものがある。具体的な実例としては、アメリカ大統領選の際の州ごとに赤と青で色分けされた地図などが挙げられる。

`folium.plugin.HeatMap` メソッドを用いることで地図上の任意の位置を中心としたヒートマップを作成することができる。引数として値を与えることで半径や色の透明度、グラデーション、ぼかしの程度が設定できるほか、前述の `LayerControl` と組み合わせることで表示の切り替えも行うことができる。

大量のマーカーをまとめて表示

前述の `folium.Marker` メソッドでは、1つのマーカーに対して多くの情報を付与する方法について解説したが、`Plugins.MarkerCluster` メソッドでは、大量のマーカーを点として地図上にプロットし、一定の閾値を定めることで地図の縮尺によってその付近にあるマーカーを1つのマーカーとして表示することができる。これによって、大量のマーカーを一度にプロットした際でも見やすく情報を提供することができる。また、このような機能はマーカーの密度という観点でヒートマップのような表し方と考えることもできる。

以上で解説した各種機能の実装例については図 3.2 にて示す。また、以上のような機能のほかにも、Folium によって実現可能な Web-GIS の機能は多くあるが、本研究では特にマーカーによる情報のフィードバックおよび重ね合わせを中心に行っていく。システムの詳細については 4 章にて示す。

提案手法

§ 4.1 多様な要因を考慮したデータセットの作成

本研究では、予測する対象地域を富山県とし、過去の犯罪発生データから、特定の日にどこで犯罪が発生するかを予測するモデルを作成する。犯罪発生予測モデルを作成する際に必要なデータセットを作成するまでの流れを図 4.1 に示す。

説明変数の選定

犯罪の発生にはさまざまな要因が考えられる。そのため、できるだけ多くの説明変数を考慮することが望ましい。しかしながら、むやみやたらに目的変数と相関がない説明変数を追加しても、予測精度が上がるどころか、計算コストが増大するだけであろう。また、本システムで統計データを取得するために利用している e-Stat は、グリッドセル・小地域ごとのデータに絞っても、200 以上公開されている。さらに、それらを別のデータと計算し、新たな説明変数を作成することを許せば、その組み合わせは考慮しきれない。

そこで、機械学習による犯罪予測モデルを作成する際に、説明変数を容易に選定することができるようになっている。選定できるデータは、e-Stat で公開されているグリッドセル・小地域ごとの統計データ、および、NAVITIME で公開されている施設データである。前者については、異なるデータ間で計算した結果を説明変数として使用できるようになっている。

異なる空間的解像度のデータの結合

e-Stat で提供されている統計データは、集計されている区分ごとに、全国ごと、都道府県ごと、市区町村ごと、”…丁目”といったの小地域ごと、グリッドセルごとの 5 種類が存在する。本システムでは、予測する空間的な単位をグリッドセルとしているため、グリッドセルごとの統計データを使用するが、より考慮できる要因を増やすため、小地域ごとの統計データも使用できるようにした。小地域ごとのデータはそのまま用いることはできないため、グリッドセル単位に変換する必要がある。そのため、小地域ごとのデータについて、小地域全体に均等に分布していると仮定し、対象のグリッドセルに重なっている割合だけを足し合わせる。すなわち、対象のグリッドセル C に重なる小地域 A_1, A_2, \dots, A_n について、それぞれの全体の面積を S_1, S_2, \dots, S_n 、対象のグリッドセルと重なる面積を s_1, s_2, \dots, s_n 、データ値を x_1, x_2, \dots, x_n とすると、対象のグリッドセル C のデータ値 X を次のように算出する。

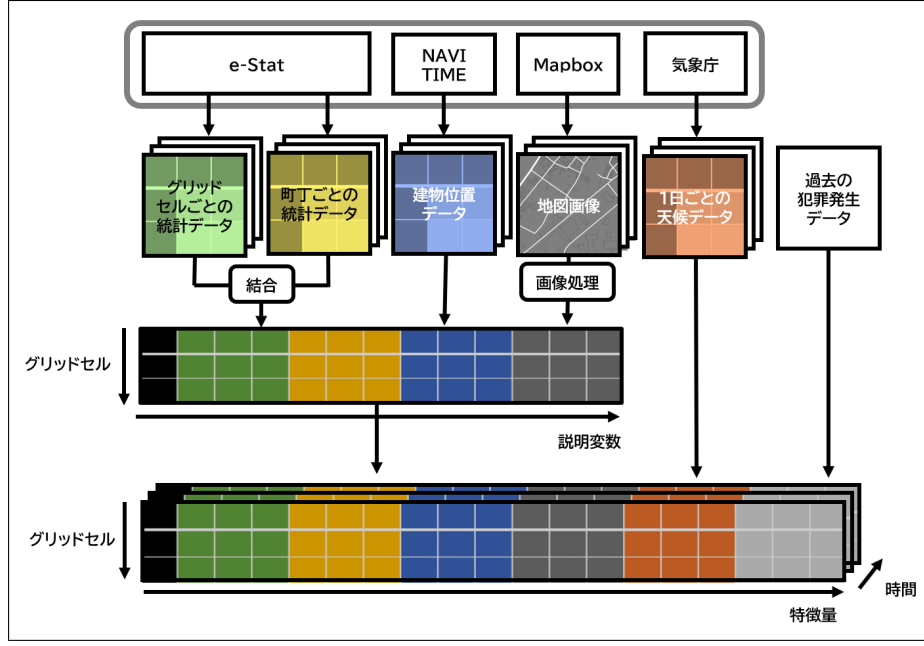


図 4.1: データセットを作成するまでの流れ

$$X = \sum_{k=1}^n x_k \frac{s_k}{S_k} \quad (4.1)$$

統計データとして公開されることの多い要素のなかには、犯罪発生の変因となり得るものが多く存在し、それらが豊富に公開されている e-Stat から、ドメイン知識をもとに自由に説明変数を選択できるようにしたことは、大きな利点と考える。

施設の最短距離と立地数の算出

さまざまなジャンルの施設について、NAVITIME からスクレイピングを行い、施設名と、その緯度と経度を取得する。本システムでは、それぞれのジャンルごとに、対象のグリッドセルに含まれる数と、最も近くにある施設までの距離を説明変数とする。

なお、地球は楕円体であるため、単純なユークリッド距離では誤差が生じてしまう。そこで、本システムではヒュベニの公式 [33] を用いて、距離を算出している。対象のグリッドセルの中心を $P_o(x_o, y_o)$ 、注目する施設を $P_n(x_n, y_n)$ とすると、それら 2 点間の距離 D は以下で求まる。

$$D = \sqrt{(D_y M)^2 + (D_x N \cos P)^2} \quad (4.2)$$

$$M = \frac{R_x(1 - E^2)}{W^3} \quad (4.3)$$

$$N = \frac{R_x}{W} \quad (4.4)$$

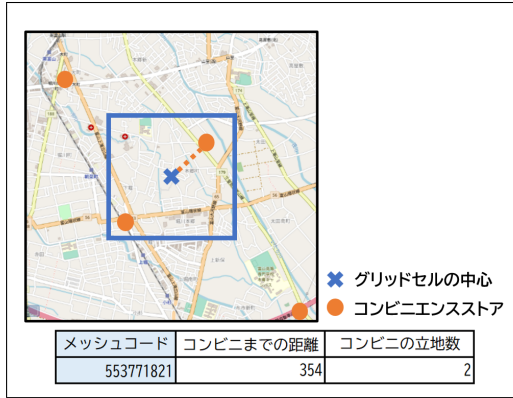


図 4.2: 施設データに基づく説明変数

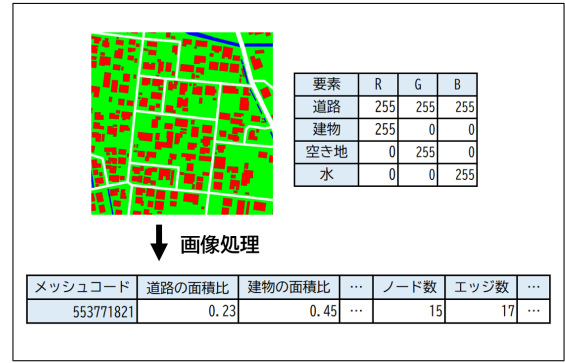


図 4.3: 地図画像に基づく説明変数

$$W = \sqrt{1 - E^2 \sin^2 P^2} \quad (4.5)$$

$$E = \sqrt{\frac{R_x^2 - R_y^2}{R_x^2}} \quad (4.6)$$

多くの人が集まりやすいレジャー施設やショッピングモールは、犯罪生成・誘引要因となりやすい。逆に、人気が少ない駐車場は、犯罪可能要因となり得る。施設に関する説明変数を自由に選択できるようにしたことは、犯罪要因の特定に役立つことが期待できる。

地図画像にもとづく説明変数の抽出

Mapbox Static Tile API を利用して、それぞれのメッシュに対応する地図画像を取得する。本研究では、建物、道路、水、空き地の4つを色分けした地図画像を取得し、それぞれの画像の大きさに対する面積の比率を説明変数としている。他人による自然な監視は犯罪を抑制する。たとえば、道路や建物の面積比率が大きいほど、監視の量が増え、犯罪が起りにくく、逆に空き地の面積比率が大きいほど、犯罪が発生しやすい傾向があるならば、それぞれの説明変数は有用なものとなるだろう。

対象の要素の面積比率 p_a は、地図画像の大きさを $n \times m$ 、その要素と同一の RGB 値をもつピクセル数を x とすると、以下のように算出する。

$$p_a = \frac{x}{nm} \quad (4.7)$$

なお、Mapbox Static Tile API によって取得する地図画像は、本来は画像処理を目的としていない。そのため、Mapbox Studio 上で指定した RGB 値と誤差があるピクセルがある。そのため、対象のピクセルの RGB 値と、それぞれの要素の RGB 値とのユークリッド距離を算出し、最も小さい要素を指定する。

また、地図画像から道路ネットワークを抽出し、道路に関連する属性を説明変数として抽出する。まず、道路とそれ以外の2値画像に変換し、ノイズを削除するためにオープニング処理を行う。その画像に対して、ネットワークを抽出するアルゴリズム [34] を使用し、ネットワークの属性であるノード数 N 、エッジ数 E を取得する。また、それらから密度 d 、平均次数 k を以下のとおり算出する。

$$d = \frac{2E}{N(N-1)} \quad (4.8)$$

$$k = \frac{2E}{N} \quad (4.9)$$

なお、密度 d と平均次数 k は、道路のネットワークとしてみたとき、それぞれ次のような特徴をもつ [35]。密度 d が大きい道路ネットワークは、道路が網目状に相互に接続された状態であり、幅員の狭い生活道路であると考えられる。また、平均次数 k が小さい道路ネットワークは、交差点の少ない直線的な道路が多いと考えられる。

動的データと静的データの組み合わせ

上で述べたものはすべて、1日ごとに変化しない静的データであった。それらと、1日ごとに変化する動的データから、空間（グリッドセル）軸 × 特徴量 × 時間軸の、3次元のテーブルを作成する。本システムでは、動的データとして、対象のグリッドセルやその周囲で、過去一定期間に発生した犯罪発生件数や、平均気温、日照時間、降水量、降雪量などの天候データを用いる。前者については、犯罪の発生には近接反復被害効果があることから採用した。また、後者については、たとえば、雨や雪が降っている日は、自転車を使う人が減少し、それと同時に自転車盗難も減少するなど、天候も少なからず犯罪発生に寄与すると考え、採用した。

以上により、機械学習に用いるデータセットの作成が完了する。

§ 4.2 不均衡なデータに対処した予測モデルの構築

作成したデータセットを用いて、犯罪発生予測モデルを作成する。犯罪が発生するデータが少なく、不均衡であることを考慮し、本システムでは XGboost を機械学習のアルゴリズムとして用いることとする。

XGBoost は、アンサンブル学習のひとつであり、2 値分類の場合は、以下のような流れで学習する。

- 1 予測値 y を求める。ただし、 $0 \leq y \leq 1$ であり、初期値は $y = 1$ である。
- 2 目的変数と予測値との誤差を最小にする決定木を構築する。
- 3 目的変数と予測値の誤差から各ノードの出力値を求める。
- 4 各ノードの出力値から予測値を計算する。
- 5 予測値と目的変数との誤差を計算する。
- 6 それを目的変数にとし、目的変数と予測値との誤差を最小にする決定木を構築する。
- 7 2～5 を繰り返し、誤差を最小化するように逐次的に学習が進む。

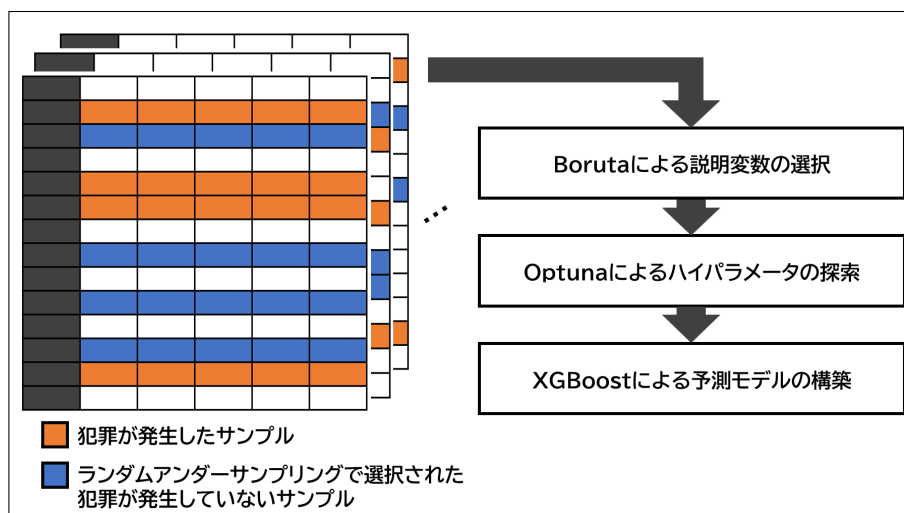


図 4.4: 予測モデルを作成する流れ

適切なランダムダウンサンプリングの探索

不均衡なデータに対するアプローチのひとつにアンダーサンプリングがある。本システムでは、犯罪が発生していないデータからランダムに抽出する、ランダムアンダーサンプリングを行う。

しかしながら、犯罪にはホットスポットと呼ばれる概念があり、特定の数少ない地域に多く発生する傾向がある。そのため単純に、犯罪が発生していないデータの数を、犯罪が発生したデータの数を同一になるようにダウンサンプリングを行った場合、予測対象の地域全体に対して、データセットに含まれる地域の割合は小さくなり、予測精度が小さくなる可能性が考えられる。

本システムで用いるデータセットは、1日単位の時間軸をもっているため、1日ごとにランダムダウンサンプリングを行い、新たなデータセットを作成する。このとき、犯罪が発生していないデータからサンプリングする数は、同日に発生したデータの数と同数にしたものとする。

Boruta による説明変数の選択

本システムでは、ユーザが自由に説明変数を選択することができるが、過度に説明変数の数が大きかったり、目的変数と相関がない説明変数があると、過学習などによって、かえって予測精度が低下してしまう可能性がある。そこで、本システムでは、犯罪発生予測モデルを作成する前に、Boruta [36] と呼ばれるアルゴリズムを用いて、適切な説明変数を選択する。

Boruta のアルゴリズムは、以下のような流れである。

- 1 もともとのテーブルをコピーし、各列をシャッフルする。もとのテーブルの説明変数を Original features, シャッフルした説明変数を Shadow features と呼ぶことにする。このとき、Shadow features は、なんら目的変数に寄与しないはずである。

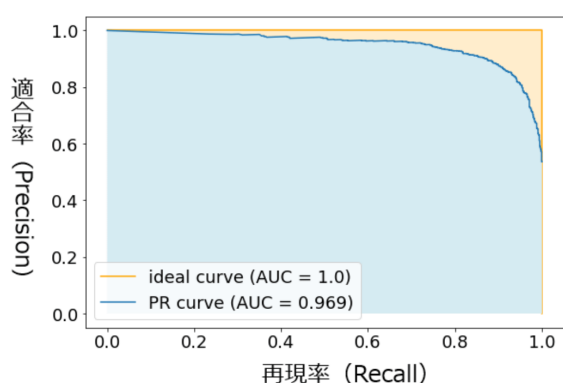


図 4.5: PR-AUC のグラフ例 [37]

		予測値	
		0	1
実測値	0	TN	FP
	1	FN	TP

$$\text{再現率} = \frac{TP}{TP+FN} \quad \text{適合率} = \frac{TP}{TP+FP}$$

$$F_1\text{スコア} = \frac{2TP}{2TP+FP+FN}$$

図 4.6: 混同行列と評価指数

- 2 Original features と Shadow features を結合し，ランダムフォレストでモデルを作成する．
- 3 そのモデルにおいて，それぞれの説明変数の重要度を算出し，Shadow features における最大値よりも大きい Original features を見つける（hit する）．
- 4 ランダムフォレストの性質により，モデルを作成するごとに重要度は変化するため，1～3 を n 回繰り返す．
- 5 各 Original features について，Shadow features の重要度と同じことを帰無仮説，より大きい・より小さいことを対立仮説とし，hit した合計を検定統計量 T ， $p = 0.5$ としたときの二項分布を用いて検定を行う．

検定の結果，説明変数が，Confirmed，Tentative，Rejected の3つに分類される．本システムでは，Confirmed，Tentative の2つを，説明変数として用いることとする．

Optuna によるハイパラメータの探索

本研究では，機械学習のアルゴリズムとして XGboost を用いる．XGBoost は，学習する際，決定木の数や最大深度，学習率などといった，事前に決定しなければならない項目（ハイパラメータ）が存在する．ハイパラメータは，一律に最適解は存在せず，データごとに最も予測精度が大きくなるものを探索する必要がある．

ハイパラメータの探索は，さまざまなアルゴリズムが提案されているが，本システムでは Tree Parzen Estimator (TPE) を用いることとし，それを利用できるフレームワーク “Opuna” を用いることにする．

TPE は，学習したモデルの評価基準をもとに，ベイズ最適化によってハイパラメータを探索するものである．モデルの評価基準には，Area Under the Precision-Recall Curve (PR-AUC) を用いる．XGBoost によって作成された2値分類を行うモデルは， $0 \leq y \leq 1$ を出力する．どれぐらい正例・負例を正しく予測できているかどうかは，しきい値をどのように決定するかどうかに左右される．PR-AUC は，適合率 (Precision) と再現率 (Recall) をしきい値ごとに算出し，適合率を縦軸，再現率を横軸としてプロットしたときの，下側にある面積である．なお，適合率と再現率は，それぞれ以下のように算出される．

$$Precision = \frac{TP}{TP + FP} \quad (4.10)$$

$$Recall = \frac{TP}{TP + FN} \quad (4.11)$$

ここで、観測値が正例のものについて、正例として予測した数を True Positive (TP)、負例として予測した数を False Negative (FN)、観測値の負例のものについて、正例として予測した数を False Positive (FP)、負例として予測した数を True Negative (TN) と表す。同様の指標として、真陽性率 (True Positive Rate) を縦軸、偽陽性率 (False Positive Rate) を横軸とする Area Under the ROC Curve (AUC) があるが、AUC-PR は正例の予測に焦点を当てた指標であるため、不均衡なデータにおけるモデルの性能差をより明確に捉えることが期待できる。

以上により、適切なダウンサンプリング、および説明変数の選択を行ったデータセットを用いて、探索によって発見したハイパラメータで XGboost による犯罪発生予測モデルを生成する。

§ 4.3 犯罪発生要因の可視化

犯罪発生予測モデルを使用し、グリッドセルごとに犯罪が発生する予測したマップと、発生する要因を示したマップを作成する。

予測精度を最大化する閾値の決定

犯罪発生予測モデルに、前日までの犯罪発生データを入力する。それらをもとに、当日に犯罪が発生するかを予測し、その結果をマップとして表示する。ここで、モデルから出力される予測値 y は $0 \leq y \leq 1$ であり、どの予測値 y から 0 または 1 とみなすか、その閾値 t を決定する必要がある。そこで、過去に発生した犯罪発生データを検証用とし、 $t = 0.001, 0.002, \dots, 0.999$ としたときの F_1 スコアを算出する。最終的に、最も F_1 スコアが大きくなったときの閾値 t を採用する。なお、 F_1 スコアは、以下のとおり算出する。

$$F_1 = 2 \frac{Precision \cdot Recall}{Precision + Recall} = \frac{2TP}{2TP + FP + FN} \quad (4.12)$$

すべてのグリッドセルに対する予測値 y に対して、採用した閾値 t が $y < t$ であれば予測値を $y = 0$ 、 $y \geq t$ であれば $y = 1$ とし、マップ上に表示する。

犯罪発生要因の可視化

犯罪発生予測モデルに対して、SHAP を適用する。SHAP は、マクロな解釈手法であり、ひとつのインスタンスについて、それぞれの説明変数の SHAP 値が算出される。犯罪発生要因マップを作成するときには、過去の犯罪発生データを入力とする。過去の犯罪発生データは、空間軸 × 特徴量 × 時間軸の 3 次元であった。そのため、これらすべてを SHAP に適

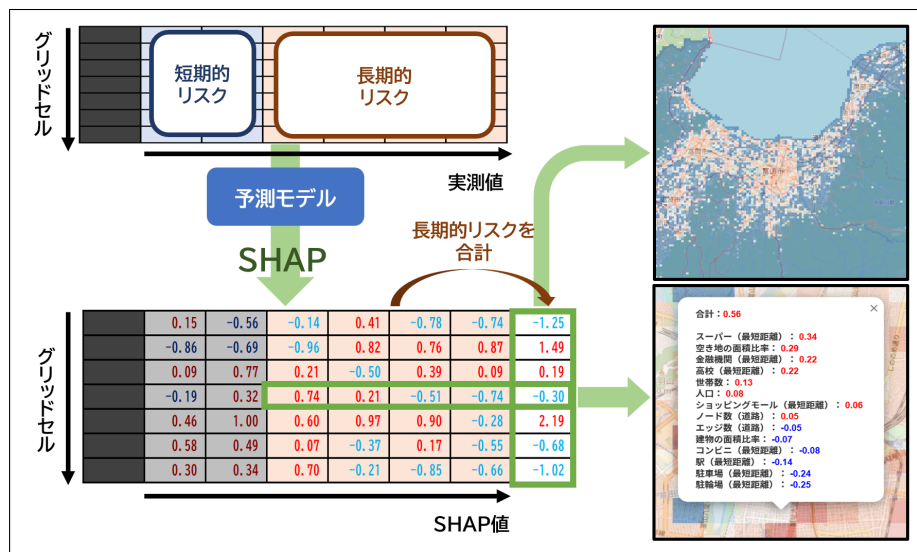


図 4.7: 要因マップを作成する流れ

表 4.1: SHAP 値テーブルの例

メッシュコード	要素1	要素2	要素3	要素4	要素5	合計
00001	0.32	0.42	-0.2	0.1	0.04	0.68
00002	-0.23	-0.1	0.23	-0.31	-0.5	-0.91

用すると、すべてのグリッドセルに対して、予測値に対する説明変数の SHAP 値が、1 日ごとに算出されることになる。

このとき、1 日ごとに変化しない静的データ、すなわち長期的リスクとなり得る要素は、SHAP 値も変化しない。そこで、出力された SHAP 値テーブルから動的データを削除し、さらに、グリッドセルごとに静的データをまとめる。これにより、それぞれのグリッドセルの長期的リスクが、犯罪発生にどれぐらい影響しているのかを表したテーブルが完成する。

SHAP 値は加法性をもっている。つまり、SHAP 値テーブルの各行の合計は、それぞれのグリッドセルが長期的リスクによって、どれだけ犯罪が発生するリスクが大きいかを表している。そして、それぞれの要素の値は、その長期的リスクがどれぐらい犯罪発生に影響しているのかを表している。

たとえば、表 4.1 を例として試してみる。まず合計の列に着目すると、上のグリッドセルは犯罪が発生しやすく、下のグリッドセルは犯罪が発生しにくいと予測していることが分かる。また、各要素の SHAP 値に着目すると、上のグリッドセルは要素 2 が最も犯罪の発生に影響しており、下のグリッドセルでは要素 4 が最も犯罪の発生を抑制する傾向があることが分かる。これをもとに、GIS 上に可視化する。図 4.8 に可視化した例を示す。SHAP 値の合計が大きくなるほどグリッドセルを赤く着色し、小さくなるほど青く着色することで、どこで犯罪が発生しやすいのかを分かりやすく可視化する。また、それぞれグリッドセルをクリックすることで、各要素の SHAP 値が大きい順で表示される。たとえば、図 4.9 の例では、スーパーとの最短距離が最も犯罪発生に影響しており、次に金融機関との最短距

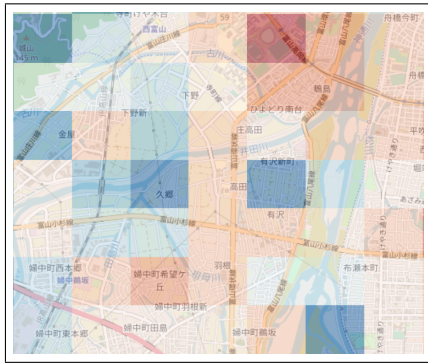


図 4.8: 犯罪発生要因マップの例

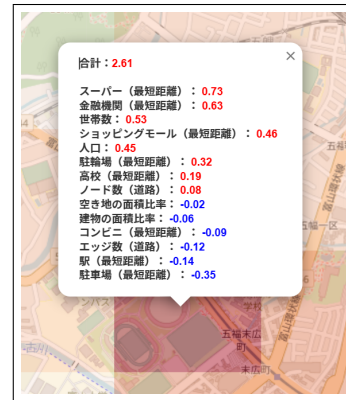


図 4.9: 各要素の SHAP 値の例

離、世帯数と続いている。逆に、駐車場との最短距離、駅との最短距離は、犯罪の発生に負の影響を与えていることが分かる。

本研究ではこのように、機械学習によって作成した犯罪発生予測モデルに対して、解釈手法の一種である SHAP を適用することによって、それぞれのグリッドセルはそれぐらい犯罪が発生しやすいのか、それぞれの要因がどれぐらい影響しているのかを GIS 上に可視化することによって、ただ単に予測の結果をもとにパトロールを強化するだけではなく、別のアプローチから犯罪を抑止することを支援する。

数値実験並びに考察

§ 5.1 数値実験の概要

本章では、実際の犯罪発生データを用いて、4章で述べた手法で犯罪発生予測モデルを作成し、その精度を確認、および考察を行う。また、その予測モデルを用いて、要因を可視化したマップを作成し、考察を行う。

犯罪発生データ

今回の数値実験で使用する犯罪発生データは、富山県警が公開している「犯罪発生マップ」[?]から取得したものをを用いる。犯罪発生データには、データ項目として「発生時刻」、「罪種」、「発生場所（緯度、経度）」が含まれており、「自転車盗難」、「ひったくり」、「車上ねらい」、「部品ねらい」、「わいせつ」、「声かけ、つきまとい」、「タイヤ盗難」の7種類の路上犯罪が収録されている。また、今回使用するレコードは、2010年9月1日から2020年9月31日の約10年間とし、発生時刻や発生場所が不正（欠損や範囲外など）なものを除外した25,814件である。なお、「犯罪発生マップ」は、このような分析を目的としておらず、掲載されている情報に誤差や欠損が存在する可能性があることに注意されたい。

予測の空間的な解像度は一辺が約500mのグリッドセル、時間的な解像度は1日とする。グリッドセルの基準としては、日本標準地域メッシュの2分の1地域メッシュを採用した。なお、富山県に含まれるグリッドセルは17,031個であるが、データセットに含まれる10年間で犯罪が発生したグリッドセルは約18.4%の3,126個であった。それ以外のグリッドセルでは、犯罪が発生する可能性は今後も限りなく小さいと考え、予測する対象のグリッドセルは、その3,126個のみとした。

データセット

予測に用いる説明変数は、表5.3に示す計65個とした。

3,126個のグリッドセルについて、静的データをまとめたテーブルを作成し、それに動的データを追加した3次元のテーブルを作成する。予測モデルの精度を検証するため、データセットのうち、2020年8月31日までの10年間を学習用、それ以降の1か月間を検証用とする（図5.1参照）。

よって、この時点におけるデータセットのレコード数は、 $N = 11419278$ であり、式??における不均衡度 r は、 $r \approx 0.0023$ である。そこで、学習用データについて、ランダムサンダーサンプリングを行い、 $N = 51065$ 、 $r \approx 0.477$ となった。

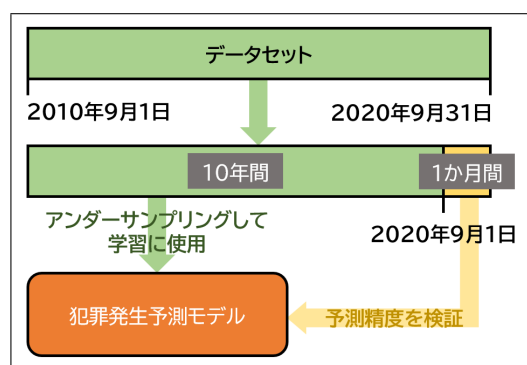


図 5.1: 予測モデルを検証する流れ

表 5.1: データセットに含まれる罪種

カテゴリ	発生件数
自転車盗難	13121
声かけ・つきまとい	6760
車上荒らし	5211
タイヤ盗難	3968
部品ねらい	690
わいせつ	401
ひったくり	59

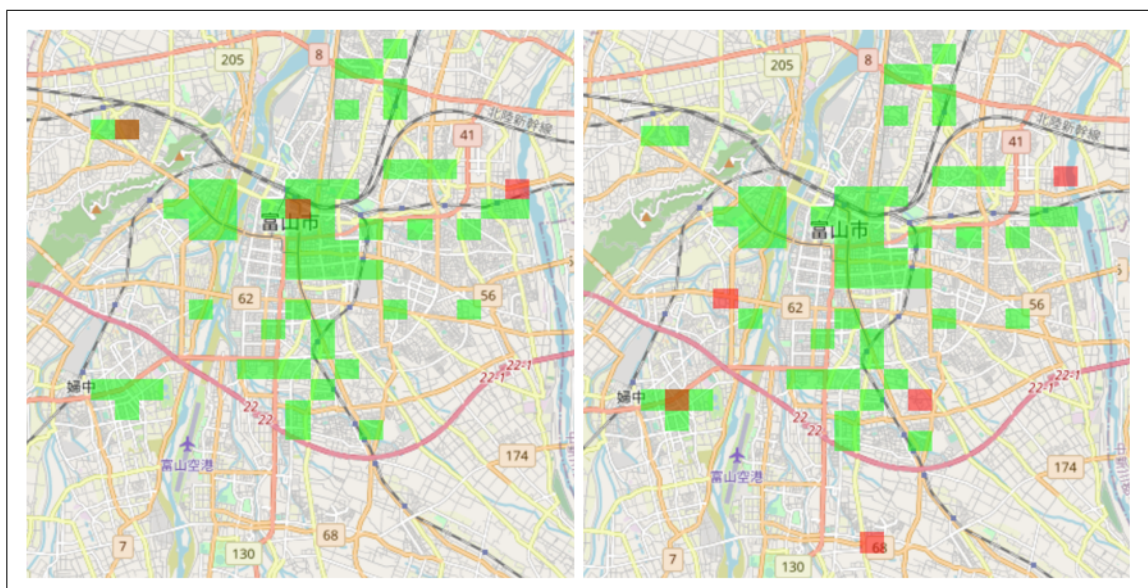


図 5.2: 2020 年 9 月 1 日（左）と 2 日（右）の予測結果

§ 5.2 実験結果と考察

予測モデルの精度検証

検証用のデータを用いて、作成した犯罪発生予測モデルの精度を検証した結果を表 5.2 に、同地域における複数日の予測結果を図 5.2 に示す。犯罪が発生したグリッドセル、発生しなかったグリッドセルともに、正しく予測した確率（正解率）は約 0.970 であったが、発生したグリッドセルを、発生すると予測した確率（再現率）は約 0.318、発生すると予測したグリッドセルで、実際に発生した確率（適合率）は約 0.018 であった。すなわち、犯罪が発生しないグリッドセルは比較的正しく予測できているものの、犯罪が発生しているグリッドセルについては、それに多少のランダム性を持っていたとしても、実用的な精度であるとは到底いえない結果となった。

この理由として、図 5.2 で分かるように、この 2 日間で犯罪が発生すると予測したグリッドセルが変化していない。表 5.3 のうち、灰色で着色した説明変数は、Boruta によって選択されたものであることを示しているが、これから分かるように、1 日ごとに変化する説明

表 5.2: 検証用データによる予測結果

		予測値	
		0	1
実測値	0	90946	2677
	1	107	50

適合率	0.018
再現率	0.318
F1スコア	0.035

変数のうち、選択されたものは「過去1か月間の犯罪発生件数」のみであった。このため、そのような静的データに対して、動的データが予測値に寄与する絶対量が小さくなり、この予測モデルは、1日ごとに予測値が変化しにくい可能性が考えられる。

短期的リスク、特に近接反復という犯罪の特性は、大きな犯罪発生 of 要因となり得る。そのため、静的データと動的データと分けて、それぞれに対して予測モデルを構築し、前者の予測モデルの予測値に対して重みづけを行うことによって、1日ごとに予測値を大きく変化させることによって、予測精度が改善する可能性がある。

また、図 5.2 に示している地域に含まれるグリッドセルは約 550 個であり、実際に犯罪が発生したグリッドセルは3~5 個である。すなわち、その割合は0.01 を下回る。本研究では、適切なアプローチを行い、不均衡なデータであっても、時空間的に解像度の大きい予測を行うことを目指したが、今回の犯罪発生データは、その限界を超えており、それでもなお精度の向上が見込めなかった可能性がある。

そのため、この改善案として、犯罪が発生する例を異常な例として、異常検知問題として取り扱うことが挙げられる。異常検知問題とは、検知したい異常な例がかなり少ないか、まったくないときに、正常な例のみを用いて、異常な例か正常な例かを判断することである。異常検知問題を取り扱うために用いるアルゴリズムは、異常な例を必要としないため、極端な不均衡、もしくは、まったく例がないデータを前提としていることである。そのため、犯罪が発生する例を異常な例として、異常検知問題として予測を行うことで、精度の向上が期待できるだろう。

犯罪発生要因の可視化

次に、学習用データを用いて、予測モデルに SHAP を適用することにより、予測モデルを可視化した結果を、図 5.3 に示す。4.3 節で述べたように、それぞれのグリッドセルに描画されている色は、各説明変数（なお、長期的リスクのみ）のもつ SHAP 値の合計を示しており、0 を基準に、大きくなるほど濃い赤色に、小さくなるほど濃い青色となっている。長期的リスクのみを抽出しているため、SHAP 値の合計は、そのグリッドセルの潜在的な犯罪発生リスクと捉えることができるだろう。

まず、学習用データを入力したときの予測モデルの精度を表 5.2 に示す。学習用データは、犯罪が発生する例としない例がほとんど同数となるようにアンダーサンプリングを行っていること、予測モデルを構築するときを使用したため、再現率は約 0.842 と大きかった。

また、図 5.3 において、赤枠で示したグリッドセルの犯罪発生要因を可視化した結果を、一例として図 5.4 に示す。これらのグリッドセルは、隣接しているのにもかかわらず、右のグリッドセルの SHAP 値の合計、すなわち犯罪が発生するリスクの合計は、左のグリッド

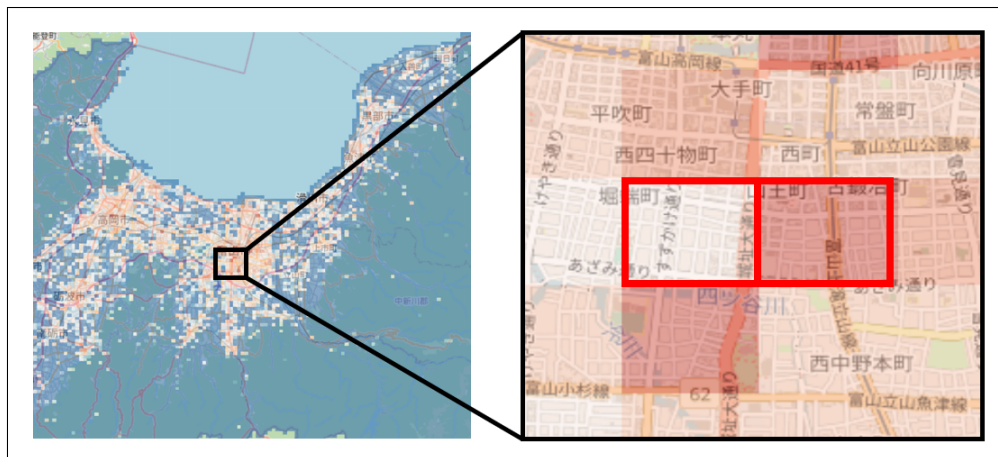


図 5.3: 予測モデルを可視化した結果



図 5.4: 特定のグリッドセルの要因を可視化した結果

セルの約 3.45 倍であると算出された。その内訳を確認すると、右のグリッドセルでは、左のグリッドセルと比較して、特に「世帯数」、「駐車場 (最短距離)」、「駐輪場 (最短距離)」の SHAP 値が増加しており、その差は、それぞれ 0.37, 0.3, 0.32 であった。すなわち、右のグリッドセルは、左のグリッドセルと比べて、世帯数が多いこと、また、駐車場、駐輪場が近い（もしくは、そのグリッドセル内にある）ことが犯罪の発生に寄与している、という知見を得ることができるだろう。



なお、本研究では、データセットとして使用した、または可視化された説明変数を「要因」と仮定し分析を行ったが、あくまでも予測値に対する説明変数の「関係性」を示しており、実際に因果関係を検証するためには、因果推論などの手法を用いる必要があることに注意する必要がある。さらに、解釈された結果は、予測モデルの精度に左右されるため、その精度が大きいことが前提となっていることにも留意すべきである。

しかしながら、単純にその場所で過去に発生した犯罪の件数を蓄積し、「ここは犯罪が発

生しやすい」と判断するだけでなく、予測値に対する各説明変数の貢献度を算出し、どのような要素が犯罪の発生に寄与しているのか、または寄与していないのか、その傾向を可視化することで、上記のような新たな知見を得られる可能性があることは、本研究における有意性のひとつと言えるだろう。一方で、可視化される要因の精度も少なからず考慮しなければならない。そこで、今後の課題として、予測モデルを介さず、実際のデータから各説明変数の貢献度を算出し、可視化することが挙げられるだろう。

表 5.3: 使用, および選択された説明変数一覧

長期的リスク		短期的リスク
人口	コンビニエンスストア(立地数)	平均気温
18歳未満人口割合	コンビニエンスストア(最短距離)	日照時間
65歳以上人口割合	駅(立地数)	降水量
外国人人口割合	駅(最短距離)	降雪量
世帯数	駐車場(立地数)	過去7日間の犯罪発生件数
単身世帯割合	駐車場(最短距離)	過去1か月間の犯罪発生件数
核家族世帯割合	駐輪場(立地数)	休日(ダミー変数)
正規労働者割合	駐輪場(最短距離)	曜日(ダミー変数)
非正規労働者割合	金融機関(立地数)	
最終学歴が中学以下の人口割合	金融機関(最短距離)	
最終学歴が高校の人口割合	旅館(立地数)	
最終学歴が大学以上の割合	旅館(最短距離)	
居住年数5年未満の人口割合	ホテル(立地数)	
居住年数20年以上の人口割合	ホテル(最短距離)	
一戸建て世帯割合	スーパーマーケット(立地数)	
アパート・低中層マンション世帯割合	スーパーマーケット(最短距離)	
高層マンション割合	ショッピングモール(立地数)	
道路の面積比率	ショッピングモール(最短距離)	
建物の面積比率	デパート(立地数)	
空き地の面積比率	デパート(最短距離)	
水の面積比率	警察署／交番(立地数)	
ノード数(道路)	警察署／交番(最短距離)	
エッジ数(道路)	小学校(立地数)	
密度(道路)	小学校(最短距離)	
平均次数(道路)	中学校(立地数)	
	中学校(最短距離)	
	高校(立地数)	
	高校(最短距離)	
	大学(立地数)	
	大学(最短距離)	
	保育園／幼稚園(立地数)	
	保育園／幼稚園(最短距離)	
	レジャー施設(立地数)	
	レジャー施設(最短距離)	

  Borutaによって選択された説明変数

おわりに

本研究では、欧米を中心に研究や実用されている地理的犯罪予測について、犯罪が発生する頻度が小さいわが国においても、適切なアプローチを行うことによって、時空間的に解像度の大きい予測を行う手法を検討した。また、予測モデルに対して、解釈手法を用いることによって、特定の地域ごとに犯罪が発生する要因を算出し、GIS上に可視化する手法を提案した。また、予測の精度をさらに高めるために、統計データなどのオープンデータのほかに、地図画像という非構造データを処理することによって、そのエリアの地理的な特徴量を抽出した。また、ナビゲーションサービスからスクレイピングをすることにより、さまざまなジャンルの施設データを取得し、犯罪発生予測モデルの構築に使用した。

数値実験では、富山県警察が公開している「犯罪発生マップ」から犯罪発生データを取得し、本研究で提案している手法を用いて、犯罪発生予測モデルを構築した。学習に使用しなかった検証用データによる予測モデルの精度の検証では、満足いく予測精度ではなかったものの、予測モデルを構築する際の特徴量選択により、地図画像から抽出した建物や空き地の面積比率、道路のエッジ数が犯罪の発生に寄与していることが分かり、地図画像からの特徴量は、予測精度に正の影響を与えることが分かった。

また、予測モデルを解釈することにより、予測モデルはどの地域で犯罪が発生しやすいと予測しているのか、また、その要因はどのようなものなのかを分かりやすく可視化した。そのため、たとえば、この地域では犯罪が発生しやすく、特に金融機関の最短距離が強く影響しているから、その地域にある金融機関を特にパトロールしたり、また金融機関に注意喚起の張り紙をつけるなど、犯罪抑止への新たな知見が得られることが期待できる。

今後の課題として、まず予測精度の向上が挙げられる。本研究の手法により作成した予測モデルは、必ずしも実用的だとは言えず、改善が必要である。たとえば、さまざまなサンプリング手法や、アンサンブル学習のアルゴリズムで比較する必要があるだろう。また、犯罪が発生することを異常だと仮定し、異常検知問題として予測することも考えられる。

また、要因の可視化は、予測モデルに基づくものであった。すなわち、その要因の精度は予測モデルの精度に左右される。そこで、予測モデルを介せず、実際のデータのみで要因を可視化する手法の開発も課題のひとつであるだろう。

さらに、警察関係者が簡単に犯罪を予測したり、要因を確認できるよう、本研究で提案した手法をバックエンドにもつシステムを作成することも、今後の課題として挙げる。

謝辞

本研究を遂行するにあたり，多大なご指導と終始懇切丁寧なご鞭撻を賜った富山県立大学工学部電子・情報工学科情報基盤工学講座の奥原浩之教授，António Oliveira Nzinga René 講師に深甚な謝意を表します．最後になりましたが，多大な協力をしていただいた研究室の同輩諸氏に感謝致します．

2023 年 2 月

島部 達哉

参考文献

- [1] 張曉齊, 米澤剛, 吉田大介, “オープンデータと LSTM を用いた犯罪発生の予測及び時間的近接性における考察”, 情報学, Vol. 16, No. 1, 2019
- [2] 国連薬物犯罪事務所 (UNODC) , ”dataUNODC”, <https://dataunodc.un.org/>, 2023 年 1 月 10 日閲覧
- [3] Cohen, L., Felson, M., and Land, K., “Property crime rates in the united states: Amacrodynamic analysis, 1947-1977; with ex ante forecasts for the mid-1980s.”, *American Journal of Sociology*, Vol. 86, No. 1, pp. 90-118, 1980.
- [4] Sherman, L., Gartin, P., and Buerger, M., “Hot spots of predatory crime: Routine activities and the criminology of place.”, *Criminology*, Vol. 27, No. 1, pp. 27-56, 1989.
- [5] Groff, E., and La Vigne, N., “Forecasting the future of predictive crime mapping”, *Crime Prevention Studies*, Vol. 13, pp. 29-58, 2002.
- [6] 雨宮護, “犯罪心理学辞典”, 丸善出版, 2016.
- [7] 大山智也, “日本における地理的犯罪予測手法の開発に関する研究”, 筑波大学システム情報工学研究科博士論文, 2020.
- [8] 野貴泰, 糸井川栄一, “犯罪多発地点の予測に基づく防犯パトロール経路に関する提案”, 地域安全学会論文集, No.31, 2017
- [9] 花岡和聖, “公然わいせつに関連する犯罪発生場所の時間的・地理的特徴：地理情報システムを活用した空間分析”, 立命館大学人文学会, Vol. 649, pp. 197-205, 2017
- [10] 花岡和聖, “大阪府における不審者遭遇情報の地理的分布: Risk Terrain Model を用いた犯罪リスクのマッピング”, 立命館大学人文学会, Vol. 656, pp. 708-720, 2018
- [11] 森本修介, 川向肇, 申吉浩, “情報伝搬モデルとガウス過程に基づく犯罪予測”, 人工知能基本問題研究会, No. 109, 2019
- [12] 西颯人, 樋野公宏, “オープンデータを用いた深層学習による犯罪発生予測の試み”, 公益社団法人日本都市計画学会 都市計画報告集, No. 16, 2017
- [13] 法務省, “平成 30 年版犯罪白書”, <https://hakusyo1.moj.go.jp/jp/65/nfm/mokuji.html>, 2023 年 1 月 10 日閲覧.
- [14] Giménez-Santana, A., Caplan, J., and Drawve, G., “Risk terrain modeling and socio-economic stratification: identifying risky places for violent crime victimization in Bogotá, Colombia”, *European Journal on Criminal Policy and Research*, Vol. 24, pp 417-431, 2018.

- [15] Bogomolov, A., Lepri, B., Staiano, J., Oliver, N., Pianesi, F., and Pentland, A., “Once upon a crime: towards crime prediction from demographics and mobile data”, *In Proceedings of the 16th international conference on multimodal interaction*, pp. 427-434.
- [16] Gorr, W., and Harries, R., “Introduction to crime forecasting. *International Journal of Forecasting*”, Vol. 19, No. 4, pp. 551-555, 2003.
- [17] Ratcliffe, J., Taylor, R., and Perenzin, A., “Predictive Modeling Combining Short and Long-Term Crime Risk Potential: Final Report”, <https://www.ncjrs.gov/pdffiles1/nij/grants/249934.pdf>, 2023 年 1 月 31 日閲覧.
- [18] Taylor, R., Ratcliffe, J., and Perenzin, A., “Can we predict long-term community crime problems? The estimation of ecological continuity to model risk heterogeneity”, *Journal of research in crime and Delinquency*, Vol. 52, No. 5, pp. 635-657.
- [19] Sampson, R., Lauritsen, J. L., “Violent victimization and offending: Individual situational- and community-level risk factors.” *Understanding and Preventing Violence*, Vol. 3, No. 1
- [20] Crowe, T., “Crime prevention through environmental design: Applications of architectural design and space management concepts.”, *Boston: Butterworth-Heinemann*, 1991.
- [21] Gorr, W., and Olligschlaeger, A., “Crime hot spot forecasting: Modeling and comparative evaluation, final project report. Washington”, DC: National Criminal Justice Reference Service, 2002.
- [22] 長瀬永遠, “証拠に基づく政策立案のためのオープンデータを利活用した Web-GIS 可視化によるデータフュージョン”, 富山県立大学学士論文, 2022.
- [23] Alberto R. Gonzales, Regina B. Schofield, Sarah V. Hart, “Mapping Crime: Understanding Hot Spots”, <https://citeseerx.ist.psu.edu/document?repid=rep1&type=pdf&doi=06b4dc6a56092493439faf0f5367293cca3f00b7#page=7>, 2023 年 1 月 31 日閲覧.
- [24] Gorr, W., and Harries, R., “Introduction to crime forecasting”, *International Journal of Forecasting*, Vol. 19, No. 4, pp.551-555.
- [25] Government Technology Magazine, “The Role of Data Analytics in Predictive Policing”, <https://www.govtech.com/data/Role-of-Data-Analytics-in-Predictive-Policing.html>, 2023 年 1 月 31 日閲覧.
- [26] Nathalie Japkowicz, “The class imbalance problem: Significance and strategies”, *In: Proc. of the Int’l Conf. on Artificial Intelligence*, 2000.
- [27] 藤原幸一, “スモールデータ解析と機械学習”, オーム社, 2022

- [28] Open Knowledge Foundation, “Place overview - Global Open Data Index”, <http://index.okfn.org/place.html>, 2023 年 2 月 23 日閲覧.
- [29] 亀谷由隆, “説明可能 AI 技術のこれまでとこれから”, 電子情報通信学会 基礎・境界サイエティ Fundamentals Review, No. 16, Vol. 2, pp. 83-92, 2022.
- [30] 森下光之助, “機械学習を解釈する技術”, 技術評論社, 2021
- [31] Scott M. Lundberg, Su-In Lee, “A Unified Approach to Interpreting Model Predictions”, *Neural Information Processing Systems*, Vol. 31, 2017
- [32] 吉田秀穂, 田嶋優樹, 今井優作, “決定木ベースモデルの解釈における SHAP 値の有用性の検討”, 人工知能学会全国大会論文集, Vol. 34, 2020
- [33] 三浦英俊, “緯度経度を用いた 3 つの距離計算方法”, オペレーションズ・リサーチ, December 2015.
- [34] “danvk/extract-raster-network: Extract a network graph (nodes and edges) from a raster image”, GitHub, <https://github.com/danvk/extract-raster-network>, 2023 年 2 月 7 日閲覧.
- [35] 向直人, “愛知県の犯罪オープンデータと地理的特徴量を利用した機械学習による犯罪種別の学習と予測”, 椙山女学園大学文化情報学部紀要, Vol. 21, pp. 109-119, 2022
- [36] Kursa, M. B., Rudnicki, W. R., “Feature Selection with the Boruta Package”, *Journal of Statistical Software*, Vol. 36, No. 11, pp. 1–13, 2010.
- [37] 一色政彦, “第 11 回 機械学習の評価関数（二値分類／多クラス分類用）を理解しよう：TensorFlow 2 + Keras (tf.keras) 入門 - @ IT”, <https://atmarkit.itmedia.co.jp/ait/articles/2103/04/news023.html>, 2023 年 1 月 31 日閲覧.

