Level Set Estimation を用いた太陽電池用シリコンのレッドゾーンの 効率的推定

Efficient estimation for red-zone in silicon wafers for solar cells using Level Set Estimation

穂積 祥太 *1	松井 孝太 * ²	沓掛 健太朗 *2	宇治原徹 * ^{3*4*2}	竹内 一郎 *1*2*5
Shota Hozumi	Kota Matsui	Kentaro Kutsukake	Toru Ujihara	Ichiro Takeuchi
*1名古屋工 Nagoya Ir	業大学 情報工学 astitute of Technology	:專攻 * ² 理化学很 RIKEN	研究所 革新知能統 Center for Advanced Int	合研究センター ^{elligence Project}
* ³ 名古屋大学 未来材料・システム研究所 Institute for Materials and Systems for Sustainability. Nagoya University				
	*4名 Department of F	古屋大学大学院二	L学研究科	, ,

*5物質・材料研究機構 情報統合型物質・材料研究拠点

Center for Materials Research by Information Integration, National Institute for Materials Science

For the task of estimating a spacial distribution of a physical quantity, it is common to fix the measurement positions to meshgrid points evenly allocated along the coordinates of the space. However, such fixed measurement positions often contain redundancy in the sense that not all the measurements in the meshgrid points are required for the target task. Especially when a measurement of the physical quantity is costly, it is thus beneficial to allocate the measurement points *adaptively* and reduce the number of measurements. In this study, we applied Level Set Estimation (LSE), which is a method to efficiently estimate the boundary position, to carrier lifetime mapping of silicon for solar cells, and estimated the low quality region. Our approach can reasonably estimate the boundary position by measuring less than 1% position compare to conventional approach.

1. 背景

太陽光電池の基板材料であるシリコンインゴッドの品質を測 る尺度の1つとして、キャリアライフタイムの値がある.キャ リアライフタイムは、光で励起されたキャリアが基底状態に 戻るまでの時間であり、太陽電池の変換効率と正の相関があ る [Sze & Ng 06].実際のインゴット製造では、材料のうちル ツボに触れている外周部分は製造過程で不純物が入ってしまい、 ライフタイムの値が低くなる.材料内でライフタイムの値が低 い範囲は特に"レッドゾーン"と呼ばれ、この範囲が狭いほど、 良い材料であるという1つの指標となっている[Ferrazza 02]. 通常、レッドゾーンの領域はインゴット表面で1点ずつライフ タイムを測定し、その結果得られるライフタイムの分布から求 められる [Hsieh et al. 14].

従来は、測定領域を格子状に区切り、各格子点上でライフタ イムを測定することでレッドゾーンの広さを評価していた.し かし、材料内の1点のライフタイム値を測定するためには相 応のコストがかかる.また、上述の測定方法ではレッドゾーン と無関係な領域にも大きな測定コストをかけてしまう可能性が あり、非常に測定効率が悪い.したがって、少ない測定回数で レッドゾーンの広さを正確に推定する方法を開発することは重 要な課題である.

本研究では,能動学習手法の1つであるレベルセット推定 (Level Set Estimation, LSE) [Gotovos 13] を用いてレッド ゾーンの境界を推定する方法を提案する.提案法では,ライフ タイム値の分布に対して統計モデルを仮定することでレッド ゾーン境界に関する統計的推測を行い,これに基づいて次に ライフタイム値を測定する点を適応的に選択する.太陽電池 インゴッドのライフタイム値データを用いた実データ実験を通 して,提案法が測定点数を劇的に減らすことができることを 示す.

2. 問題設定

あるブラックボックス関数 $f: \mathcal{X} \to \mathbb{R}$ と観測対象の候補点 $\{x_i\}_{i=1}^N$ が与えられているとする.ここで、 \mathcal{X} は測定点の空間 を表し、測定値は実数値であるとする.また、入力 $x \in \mathcal{X}$ に 対する関数値 f(x) の観測値 y は、誤差 $s \sim \mathcal{N}(0, \sigma^2)$ を伴っ て y = f(x) + s と表されるとする.本研究で考察するレベル セット推定 (Level Set Estimation, LSE)[Gotovos 13] は、あ る閾値 $h \in \mathbb{R}$ に対して、

 $\boldsymbol{x}_i \in \{\boldsymbol{x} \in \mathcal{X} \mid f(\boldsymbol{x}) > h\}$ or $\boldsymbol{x}_i \in \{\boldsymbol{x} \in \mathcal{X} \mid f(\boldsymbol{x}) \le h\}$ (1)

を判定する判別問題として定式化される.

特に、本研究で取り扱うレッドゾーン推定問題では、fとして試料表面上の2次元座標を入力、その座標におけるライフタイムの値を出力とするような2入力1出力の関数を考える。通常の回帰分析の文脈では、入力xと誤差を伴った出力yの複数のペア { (x_i, y_i) } $_{i=1}^{i=1}$ を学習データとしてfを推定することで、新たな入力に対する判別を行うことができる。しかし、材料科学をはじめとする基礎科学分野では、入力xに対する出力yを観測するのに大きなコストがかかる場合が少なくなく、本研究で対象とするライフタイム値の測定も同様である。したがって、できるだけ少ないyの観測回数(すなわち、少ないデータ数)で(1)の判別問題を解くことがしばしば要求される。

連絡先: 穂積祥太,名古屋工業大学大学院 竹内研究 室所属,466-8555 愛知県名古屋市昭和区御器所町, hozumi.s.mllab.nit@gmail.com

Algorithm 1 Level Set Estimation

Require: sample set *D*, GP prior($m(\mathbf{x}) = a, k, \sigma_0$) threshold ratio ω , accuracy parameter $\varepsilon = a/5$, max iteration T**Ensure:** predicted sets \hat{H} , \hat{L} , U_t 1: $H_0, L_0, U_0, Z_0 \leftarrow \emptyset, C_0(\boldsymbol{x}) \leftarrow \mathbb{R}, t \leftarrow 1$ 2: while $U_{t-1} \neq \phi$ and $t \leq T$ do $H_t \leftarrow H_{t-1}, L_t \leftarrow L_{t-1}, U_t \leftarrow U_{t-1}, Z_t \leftarrow Z_{t-1}$ 3: for all $\boldsymbol{x} \in D$ do 4: $h_t^{opt} \leftarrow \omega \max_{\boldsymbol{x} \in Z_{t-1}} \max(Q_t(\boldsymbol{x}))$ 5: $f_t^{pes} \leftarrow \max_{\boldsymbol{x} \in Z_{t-1}} \min(Q_t(\boldsymbol{x}))$ 6: $h_t^{pes} \leftarrow \omega f^{pes}$ 7: if $\min(Q_t(\boldsymbol{x})) + \varepsilon \ge h_t^{opt}$ then 8: $U_t \leftarrow U_t \setminus \{x\}$ 9: if $\max(Q_t) < f_t^{pes}$ then $H_t \leftarrow H_t \cup \{x\}$ 10: else $M_t^H \leftarrow M_t^H \cup \{x\}$ 11: else if $\max(Q_t(\boldsymbol{x})) - \varepsilon \ge h_t^{pes}$ then 12: $U_t \leftarrow U_t \setminus \{x\}$ 13:if $\max(Q_t) < f_t^{pes}$ then $L_t \leftarrow L_t \cup \{x\}$ 14: else $M_t^{\hat{L}} \leftarrow M_t^L \cup \{x\}$ 15:end if 16:end for 17: $Z_t \leftarrow U_t \cup M_t^H \cup M_t^L$ 18: $\boldsymbol{x} \leftarrow \arg \max(w_t(\boldsymbol{x}))$ 19: $\bar{x} \in Z_t$ 20: $y_t \leftarrow f(\boldsymbol{x}) + n_t$ Compute $\mu_t(\boldsymbol{x})$ and $\sigma_t(\boldsymbol{x})$, for all $\boldsymbol{x} \in D$ 21: $t \leftarrow t + 1$ 22. 23: end while 24: $\hat{H} \leftarrow H_{t-1} \cup M_{t-1}^H, \hat{L} \leftarrow L_{t-1} \cup M_{t-1}^L$

本研究では、f を、ブラックボックスかつ少データゆえの不 確実性の下で適応的に推定するために統計モデルを導入し、モ デルに基づいて得られる予測を用いて(1)を効率的に解くア プローチを提案する.

提案手法 3.

以下では,まずブラックボックス関数としての f の 不確実性を考慮に入れたガウス過程によるモデリン グ [Rasmussen & Williams 06] を説明し、その後 (1) の判別 問題のための提案アルゴリズムについて説明する.

3.1 ガウス過程によるモデリング

まず,ブラックボックス関数としてのfの不確実性を考慮 するため、fの統計モデルとしてガウス過程事前分布

$$f(\boldsymbol{x}) \sim GP(m(\boldsymbol{x}), k(\boldsymbol{x}, \boldsymbol{x}'))$$

を仮定する.ここで、 $GP(m(\boldsymbol{x}, k(\boldsymbol{x}, \boldsymbol{x}')))$ は平均関数 m と カーネル関数 k で特徴付けられるガウス過程を表す.本研究 では, k として Matern3/2 カーネル

$$k(\boldsymbol{x}_i, \boldsymbol{x}_j | \boldsymbol{\theta}) = \beta_f^2 (1 + \frac{\sqrt{3}r}{\beta_l}) \exp(-\frac{\sqrt{3}r}{\beta_l})$$

を用いる.ここで、 $r = \sqrt{(\boldsymbol{x}_i - \boldsymbol{x}_j)^T (\boldsymbol{x}_i - \boldsymbol{x}_j)}$ とおいた. Matern カーネルは、ガウスカーネルと比較してより非滑ら



図 1: 左図:データ1のライフタイム値のプロット,右図:デー タ2のライフタイム値のプロット.青い領域がライフタイム値 が低い低品質領域 (レッドゾーン)を表す、実験では、上図を 正解とし,提案法がどの程度復元できたかを評価指標の1つ としている.

かな関数を表現しやすいという特徴を持つ.このとき、観 測点 $x_{1:n} = \{x_1, ..., x_n\}$ に対して, 関数値ベクトル f = $(f(x_1), ..., f(x_n))$ は、定義からn変量の正規分布に従う

 $\boldsymbol{f}|x_{1:n} \sim \mathcal{N}(\boldsymbol{m}, \boldsymbol{K})$

ここで, $m = (m(x_1, ..., m(x_n))), K_{i,j} = k(x_i, x_j)$ である. 以上より、データ $\mathcal{D}_n = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$ が観測された下での新 たな入力 \boldsymbol{x} の観測値 $f(\boldsymbol{x})$ の予測分布は $\mathcal{N}(\mu_n(\boldsymbol{x}), \sigma_n^2(\boldsymbol{x}))$ と 書ける. ただし,

$$\mu_n(\boldsymbol{x}) = m(\boldsymbol{x}) + \boldsymbol{k}(\boldsymbol{x})^T (\boldsymbol{K} + \sigma^2 \boldsymbol{I})^{-1} (\boldsymbol{y} - \boldsymbol{m})$$

$$\sigma_n(\boldsymbol{x}) = k(\boldsymbol{x}, \boldsymbol{x}) + \boldsymbol{k}(\boldsymbol{x})^T (\boldsymbol{K} + \sigma^2 \boldsymbol{I})^{-1} \boldsymbol{k}(\boldsymbol{x})$$

はそれぞれ予測平均と予測分散を表す.

レベルセット推定 (LSE) 3.2

LSE では, 前節のモデリングに基づいて confidence region

$$Q_n(\mathbf{x}) := [\mu_{n-1}(\mathbf{x}) \pm \beta_t^{1/2} \sigma_{n-1}(\mathbf{x})]$$

を構成し、閾値 $h \ge Q_n(x)$ の関係から x が h を超えるかどう かの判定を行う.具体的には,閾値に対するマージンεを導入 して, $\min(Q_n) + \varepsilon \ge h$ であれば上側領域, $\max(Q_n) - \varepsilon \ge h$ であれば下側領域であると判別する.ただし、実問題では h が 未知であったり、一度どちらかの領域に判別された点でも、不 確実性が大きい場合は後の反復で領域が反転したりするといっ た状況が起こりうるため、より詳細な判別基準を設ける必要が ある.以上の考察に基づき,提案する LSE のアルゴリズムを Algorithm 1 に示す.

提案法では、未判定点の集合 U, $x_i \in \{x \in \mathcal{X} \mid f(x) > h\}$ かつ不確実性が十分小さいと判定された点の集合 $H, x_i \in$ $\{x \in \mathcal{X} \mid f(x) > h\}$ かつ不確実性が大きいと判定された点の 集合 M^H , $x_i \in \{x \in \mathcal{X} \mid f(x) \leq h\}$ かつ不確実性が十分小 さいと判定された点の集合 *L*, $x_i \in \{x \in \mathcal{X} \mid f(x) \leq h\}$ か つ不確実性が大きいと判定された点の集合 M^L の 5 つの集合 を繰り返し更新し、全ての観測候補点が H または L に振り分 けられた時点で終了する.

提案法では、最も不確実性の高い点で次の関数値を測定する (uncertainty sampling) アプローチを採用する. これは,獲得 関数

$$egin{aligned} w_t(oldsymbol{x}) &= \max(Q_t(oldsymbol{x})) - \min(Q_t(oldsymbol{x})) \ &= 2eta_n^{1/2}\sigma_{n-1}(oldsymbol{x}) \end{aligned}$$



図 2: データ1 に対する提案法の挙動. 左上段から左下段に向 けて,獲得点数が10,20,30,40 時点,右上段から右下段に 向けて,獲得点数が50,60,70,80 時点の推定レッドゾーン (水色領域)を表し. 濃青領域は,不確実性が残っていることを 表す.また,黒点はライフタイムの観測を行った座標を示す.

の最大値を与える点として特徴付けられる.アルゴリズムの各 ステップで,この基準で選ばれた *x* に対して *y* の値を観測し, ガウス過程を更新する.

4. 計算機実験

4.1 実験設定

本研究では、太陽電池インゴッドのライフタイム値を測定 した2種類のデータを解析した.データ1については試料表 面の2次元平面上に格子点を19481点、データ2については 14641点とり、その上でライフタイム値の測定を実施したもの である.計算機実験では、レッドゾーンの閾値をライフタイム 値の最大値の15%値と定め、これを下回る範囲をレッドゾー ンと定義する.図1にそれぞれのデータにおけるライフタイ ム分布のプロットを示す.図の青い領域が真のレッドゾーンを 表している.本実験ではライフタイムを評価する点の数の最大 値を100に固定し、アルゴリズムが停止した時点で推定され



図 3: 上段:F 値のプロット,下段:推定されたレッドゾーンの 可視化 (水色領域). 左側がデータ 1,右側がデータ 2 の結果 をそれぞれ示す.

たレッドゾーンを F 値及び可視化によって評価した.

ガウス過程に基づく LSE では,観測誤差の分散 σ² をデー タから推定するとともに,次の 5 種類のハイパーパラメータ を設定する必要がある:

- 1. カーネル関数の分散 β_f
- 2. カーネル関数の長さスケール β_l
- 3. ガウス過程の平均関数 m(x)
- 4. LSE における分類マージン ε
- 5. confidence region Q_n で平均と標準偏差のバランスを制 御するパラメータ β_t .

文献 [Gotovos 13] では, ライフタイム値の取りうる範囲を被 覆するように取ったデータの一部をを用いて, 対数尤度最大 化に基づいてパラメータの設定を行っている. しかし, 実際 の実験ではライフタイム値の取りうる範囲が未知であるため, この方法でパラメータを調整しておくことが難しい. 本研究 では, 経験的な知見から $\sigma = 0.001, \beta_l = 25, \beta_t = 1.96$ と設 定した. 他のパラメータ設定は, まず一様ランダムに 10 点の ライフタイム値を測定し, その中央値の半分の値を a とおき, $\beta_f^2 \in [25a, 100a], m(x) = a, \varepsilon = a/5$ とするヒューリスティク スを用いた. 特に β_f に関しては, 各反復で 1 点のライフタイ ムを測定する毎に上記の範囲内で対数尤度最大化に基づいて適 切な値を設定した.

4.2 実験結果

計算機実験は1000回の試行を行い、平均的な振る舞いを観察した.まず、データ1に対する提案法の挙動を図2に示す. 探索初期では大部分が濃青の不確実な領域を占めているが、探索が進むにつれて不確実性が減少していき、徐々にレッドゾーンが推定されていっていることが見てとれる.また、実際にライフタイム値を観測する点も適応的に選択されており、格子点上を全探索する従来のアプローチに比べて探索が効率的になっ ていることがわかる.次に,データ1及びデータ2に対する 結果をそれぞれ図4.に示す.図4.上段のF値のグラフは中 央値を中心にエラーバーの上限は第一四分位数,下限は第三四 分位数として描画した.F値は,データ1で0.7325,データ 2で0.8458程度を達成した.また,推定されたレッドゾーン を可視化した図4.の下段と真のレッドゾーンを示した図1及 び図4.の下段を比較すると,最大100点の測定で比較的外形 が良く推定できていることがわかる.ここで,図4.下段の濃 青の領域は,境界周辺の不確実性がまだ残っていることを表し ている.真のレッドゾーンを同定するために測定したデータ点 数に対しては,データ1で0.5% (100/19481),データ2で 0.6% (100/14641)の測定点数で上記の推定結果を得ることが できた.探索コストを大幅に削減できていることがわかる.

5. まとめ

本研究では、太陽電池インゴッドの品質管理において、 ライ フタイム値の低い低品質領域 (レッドゾーン) を推定するため に、レベルセット推定に基づく能動学習を提案した. レベルセッ ト推定は、ライフタイムをブラックボックス関数としてガウス 過程でモデリングすることで、レッドゾーンの推定とデータの 測定を適応的に行うことができるため,探索効率の改善が期待 できる.計算機実験によって、レッドゾーンを正確に同定でき る一方で探索コストの非常に高い格子点上測定に対して,提案 法では非常に少ない探索コストでレッドゾーンの概形を推定す ることができた.本研究における課題として、レベルセット推 定は欠損値を含むデータには適用できない、という点が挙げら れる.今回の実験では、データに対して専門家の知識を利用し た欠損値処理などの前処理が多数行われているが,実際の実験 でリアルタイムにレベルセット推定を行う場合、これらの前処 理を同時に実施することは難しい. 今後は、 欠損値を含むデー タに対して適用できるようにレベルセット推定アルゴリズムを 拡張し、リアルタイムの実験に適用したい.

参考文献

- [Ferrazza 02] Ferrazza, F. (2002). Large size multicrystalline silicon ingots. Solar energy materials and solar cells, 72(1-4), 77-81.
- [Gotovos 13] Gotovos, A., Casati, N., Hitz, G., & Krause, A. (2013). Active learning for level set estimation. In IJCAI (pp. 1344-1350).
- [Hsieh et al. 14] Hsieh, C. C., Lan, A., Hsu, C., & Lan, C. W. (2014). Improvement of multi-crystalline silicon ingot growth by using diffusion barriers. Journal of Crystal Growth, 401, 727-731.
- [Rasmussen & Williams 06] Rasmussen, C. E., & Williams, C. K. (2006). Gaussian process for machine learning. MIT press.
- [Sze & Ng 06] Sze, S. M., & Ng, K. K. (2006). Physics of semiconductor devices. John wiley & sons.