

Newton-CG 法による条件付き確率場のバッチ学習

坪井 祐太 海野 裕也
日本アイ・ビー・エム株式会社
{yutat,yunno}@jp.ibm.com

鹿島 久嗣
東京大学
kashima@mist.i.u-tokyo.ac.jp

岡崎 直観
東京大学
okazaki@is.s.u-tokyo.ac.jp

1 導入

直鎖条件付き確率場 (CRF) は出力ラベル列の予測器として、自然言語処理 (NLP) において広く応用されている [4]. CRF は観測 $\mathbf{x} \in \mathbf{X}$ に対するラベル列 $\mathbf{y} \in \mathbf{Y}$ の条件付き確率を $P_{\theta}(\mathbf{y}|\mathbf{x}) = e^{\theta^{\top} \Phi(\mathbf{x}, \mathbf{y})} / Z$ でモデル化する. ただし, $\theta \in \mathbb{R}^d$ は d 次元のパラメータ・ベクトル, $\Phi(\mathbf{x}, \mathbf{y}) : \mathbf{X} \times \mathbf{Y} \rightarrow \mathbb{R}^d$ は素性関数, 分母 $Z = \sum_{\mathbf{y} \in \mathbf{Y}} e^{\theta^{\top} \Phi(\mathbf{x}, \mathbf{y})}$ は分配関数である. また, \mathbf{v}^{\top} は \mathbf{v} の転置である. 訓練データ D が与えられたとき, CRF のバッチ学習は罰則項付き負の対数尤度関数 $f(\theta)$ の最小化問題として与えられる.

$$f(\theta) = -\sum_{(\mathbf{x}, \mathbf{y}) \in D} [\theta^{\top} \Phi(\mathbf{x}, \mathbf{y}) - \ln Z] + \frac{\|\theta\|^2}{2\sigma^2} \quad (1)$$

ただし, 最終項は過学習を防ぐための平均 0 分散, σ^2 の θ の正規事前分布である.

非線形関数 $f(\theta)$ の最適解は解析的には求められないため, 反復法によって最適化する. 反復法では, 各反復 k において現在のパラメータ θ_k をある探索方向 \mathbf{s}_k に沿って次の値 $\theta^{k+1} = \theta_k + \mathbf{s}_k$ に更新することを繰り返し, 最適解に収束させる. もっとも明らかな探索方向としては勾配方向 $-\mathbf{g}_k \equiv -\partial_{\theta} f_k$ がある (最急降下法). ただし, $f_k \equiv f(\theta_k)$ と略記する. また, 別の探索方向としては, Newton 方向 $-\mathbf{H}_k^{-1} \mathbf{g}_k \in \mathbb{R}^d$ がある. ただし $\mathbf{H}_k \equiv \partial_{\theta^2} f_k$ はヘシアン行列である. Newton 方向は f_k の 2 次のテイラー近似: $f(\theta_k + \mathbf{s}) \approx f_k + \mathbf{g}_k^{\top} \mathbf{s} + \frac{1}{2} \mathbf{s}^{\top} \mathbf{H}_k \mathbf{s}$ を最小化する更新 \mathbf{s} である. 2 次近似による最適化手法は非常に少ない反復数で収束することが知られている.

NLP では非常に高次元のパラメータを扱う必要がある, $d \times d$ 行列のヘシアン逆行列 \mathbf{H}^{-1} を陽に計算することは現実的でない. そこで, CRF のバッチ学習には記憶制限 Newton 法 (LBFGS) が広く使われている [8]. LBFGS は, \mathbf{H}^{-1} を直近 $m (\ll d)$ のパラメータと勾配の変化履歴 (d 次元ベクトル $2m$ 個) で近似し, Newton 方向を求める.

一方, Newton 方向は連立 1 次方程式 $\mathbf{H} \mathbf{s} = -\mathbf{g}$

(Newton 方程式) の解でもある. 本稿では Newton 方程式を共役勾配 (CG) 法によって解くことで Newton 方向を求める Newton-CG 法 [7] を CRF のバッチ学習に応用することを提案する. なお, 本手法でも \mathbf{H}^{-1} を陽に持つ必要はない. Newton-CG 法はロジスティック回帰の学習において LBFGS に勝ることが示されており [5], ロジスティック回帰の一般化である CRF においてもその有効性が期待できる. 本研究の貢献は, CRF の a) ヘシアン・ベクトル積 $\mathbf{H} \mathbf{s}$ の動的計画法による計算手順と, b) 勾配計算の中間結果の再利用による Newton-CG 法の高速度化である. 基本句チャンキングと固有表現抽出タスクでの実験において, LBFGS より数倍早く収束することが示された. さらに, 小・中規模のタスクにおいてはオンライン学習手法よりも早く最良の性能を得られることが確認された.

2 Newton-CG 法

Newton-CG 法は 2 重ループ, a) 現在のパラメータ θ_k での Newton 方向を CG 法により求める内側のループと, b) Newton 方向に従ってパラメータを θ_{k+1} に更新する外側のループで構成されている. 以降では, 内側ループの反復を ℓ , 外側ループの反復を k で示す.

内側ループで解く必要のある方程式 $\mathbf{H} \mathbf{d} = -\mathbf{g}$ の解 (Newton 方向) は, 2 次形式 $q(\mathbf{d}) = \frac{1}{2} \mathbf{d}^{\top} \mathbf{H} \mathbf{d} + \mathbf{g}^{\top} \mathbf{d}$ を最小化する解と等しい. ただし, 内側ループでは θ_k は固定であり, 勾配ベクトルとヘシアン行列は変わらないため $\mathbf{g} \equiv \mathbf{g}_k$, $\mathbf{H} \equiv \mathbf{H}_k$ と略記する. CG 法はすべての $\ell \neq h$ において \mathbf{H} に関して直交 ($\mathbf{d}_{\ell}^{\top} \mathbf{H} \mathbf{d}_h = 0$) する点列 $\mathbf{d}_{\ell} (\ell = 1, 2, \dots)$ により最悪 d 回の反復数で 1 次方程式を解くことができる. 多くの場合 d より少ない反復数で収束することが知られており, $\partial_{\mathbf{d}} q$ のノルム $r_{\ell} = \|\mathbf{H} \mathbf{d}_{\ell} + \mathbf{g}\|$ が停止条件として使われる. また, f のような非線形関数では Newton 方向は最適解近くでのみ早い収束を保障し, 初期は厳密に Newton 方向を求める必要がないため, Newton-CG 法では $\mathbf{d}_0 = -\mathbf{g}_k$ を初期点として $r_{\ell} \leq \xi_k \|\mathbf{g}_k\|$ を停止条件とする. ただ

し $0 < \xi_k < 1$ かつ $\lim_{k \rightarrow \infty} \xi_k = 0$. つまり、最適解に近づく（勾配が小さくなる）につれて正確に Newton 方向を求める.

外側ループでは、非線形関数 f の最適解に収束することを保障するために、線形探索法や信頼区間法などにより Newton 方向を調整する. 4 節の実験では信頼区間法を実装したが、紙面の制限から、Newton-CG 法のアルゴリズム詳細については最適化法の教科書 [7] 等を参照されたい.

Newton-CG 法の計算時間は CG 法が大部分を占めるが、先に述べた適応的な CG 反復の停止条件により初期の反復では勾配方向を使い、反復が進むにつれて厳密に Newton 方向を求めるため計算効率が良い. また、CG 法ではヘシアン行列 \mathbf{H} を陽に求める代わりに、 \mathbf{H} とベクトル \mathbf{d}_ℓ の積 $\mathbf{H}\mathbf{d}_\ell \in \mathbb{R}^d$ ベクトルが求められれば q の最小化が行えるため空間効率も良い. 次節では、CRF 学習の目的関数 (1) のヘシアン・ベクトル積 $\mathbf{H}\mathbf{d}$ の計算法を提案する.

3 CRF のヘシアン・ベクトル積

直鎖 CRF は κ 次マルコフ性を仮定することで、 $O(T|Y|^{\kappa+1})$ 時間でヘシアン・ベクトル積を計算できることを示す. ただし、 T はラベル列 \mathbf{y} の長さである. なお、CRF の目的関数の勾配 \mathbf{g} も同じく $O(T|Y|^{\kappa+1})$ で計算できることが知られている [4]. ヘシアン・ベクトル積は勾配計算時に算出される周辺確率を利用して高速に計算できることも示す.

任意のベクトル \mathbf{d} と目的関数のヘシアン行列の積は、事例毎の線形和として次のように分解できる.

$$\mathbf{H}\mathbf{d} = \left[\sum_{(\mathbf{x}, \mathbf{y}) \in D} \mathbf{H}(\mathbf{x}, \mathbf{y})\mathbf{d} \right] + \frac{\mathbf{d}}{\sigma^2}$$

また、事例ごとのヘシアン・ベクトル積は次式のとおりのである.

$$\mathbf{H}(\mathbf{x}, \mathbf{y})\mathbf{d} = \sum_{\tilde{\mathbf{y}} \in \mathcal{Y}} P_\theta(\tilde{\mathbf{y}}|\mathbf{x}) \Phi(\mathbf{x}, \tilde{\mathbf{y}}) \Psi_\theta(\mathbf{x}, \tilde{\mathbf{y}})^\top \mathbf{d} \quad (2)$$

ただし、 $\Psi_\theta(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x}, \mathbf{y}) - \mathbb{E}_{P_\theta}(\Phi(\mathbf{x}, \mathbf{Y}))$.

式 (2) は指数個の要素の集合 \mathcal{Y} の和を含むが、CRF の勾配計算と同様に動的計画法で多項式時間で計算可能であることを次に示す. 説明のため 1 次のマルコフ性を仮定する. マルコフ性を仮定すると素性関数は $\Phi(\mathbf{x}, \mathbf{y}) = \sum_{t=1}^{T+1} \phi(\mathbf{x}, y_{t-1}, y_t)$ と分解でき、条件付き確率は次のように分解できる.

$$P(\mathbf{y}|\mathbf{x}) = P(y_1|y_2, \mathbf{x}) \cdots P(y_{t-2}|y_{t-1}, \mathbf{x}) \\ P(y_{t-1}, y_t|\mathbf{x}) P(y_{t+1}|y_t, \mathbf{x}) \cdots P(y_T|y_{T-1}, \mathbf{x})$$

ただし、列頭と列末をそれぞれ S と E であらわすとき、 $P(S|y_1, \mathbf{x}) = 1$ と $P(E|y_T, \mathbf{x}) = 1$ は省略している. これらの分解により、式 (2) は次式で書き直すことができる:

$$\mathbf{H}(\mathbf{x}, \mathbf{y})\mathbf{d} = \sum_{t=1}^T \sum_{i, j \in Y} \phi(\mathbf{x}, i, j) M(t, i, j).$$

ただし、 $M(t, i, j)$ は次式で定義される.

$$M(t, i, j) = P_\theta(y_{t-1}=i, y_t=j|\mathbf{x}) [\mathbf{A}[t-1, i] + \phi(\mathbf{x}, i, j)^\top \mathbf{d} + \mathbf{B}[t, j] - \mathbb{E}_{P_\theta}(\Phi(\mathbf{x}, \tilde{\mathbf{y}}))^\top \mathbf{d}].$$

また、行列 $\mathbf{A}[t, j]$ と $\mathbf{B}[t, i]$ は次の再帰式で定義される.

$$\mathbf{A}[t, j] = \begin{cases} 0 & \text{if } t=0 \\ \phi(\mathbf{x}, S, j)^\top \mathbf{d} & \text{else if } t=1 \\ \sum_{i \in Y} P_\theta(y_{t-1}=i|y_t=j, \mathbf{x}) (\phi(\mathbf{x}, i, j)^\top \mathbf{d} + \mathbf{A}[t-1, i]) & \text{otherwise,} \end{cases}$$

$$\mathbf{B}[t, i] = \begin{cases} \phi(\mathbf{x}, i, E)^\top \mathbf{d} & \text{if } t=T \\ \sum_{j \in Y} P_\theta(y_{t+1}=j|y_t=i, \mathbf{x}) (\phi(\mathbf{x}, i, j)^\top \mathbf{d} + \mathbf{B}[t+1, j]) & \text{otherwise.} \end{cases}$$

よって、 \mathbf{A} と \mathbf{B} を予め計算しておくことで、式 (2) は $O(T|Y|^2)$ で計算可能である. なお、 \mathbf{A} と \mathbf{B} の式に現れる条件付き確率は、周辺確率の割り算

$$P_\theta(y_{t-1}|y_t, \mathbf{x}) = \frac{P_\theta(y_{t-1}, y_t|\mathbf{x})}{P_\theta(y_t|\mathbf{x})}$$

で計算できる. $P_\theta(y_{t+1}|y_t, \mathbf{x})$ も同様. また、期待値 $\mathbb{E}_{P_\theta}(\Phi(\mathbf{x}, \mathbf{Y}))$ も周辺確率を用いて計算可能である.

重要な点として、この周辺確率は勾配計算時に算出しているため、 θ が固定であれば異なる \mathbf{d} でのヘシアン・ベクトル積の計算で共用できる. CRF の周辺確率の計算には四則演算より数十倍遅い指数計算を含んでいる [2]. よって、周辺確率を再利用することにより、四則演算だけで計算できるヘシアン・ベクトル積は勾配計算にわずかな計算を追加するだけで計算できる. 2 節で指摘したように Newton-CG 法の内側ループ (CG 法) では θ は固定であるため、周辺確率を保存して再利用することができ、相性が良い. Newton-CG 法は内側ループの計算量が大部分を占めるため、CRF 全体の学習全体の高速化が期待できる. なお、周辺確率は各事例 (\mathbf{x}, \mathbf{y}) について κ 次 CRF では $O(T|Y|^{\kappa+1})$ の保存空間が必要になる. 利用可能なメモリ領域に合わせて訓練データ D の部分集合のみの周辺確率をキャッシュすることも可能である.

タスク	訓練	開発	テスト	$ X $	$ Y $
チャンキング	8,936	N/A	2,012	338,539	23
固有表現抽出	8,322	1,914	1,516	99,135	9

表 1: データセットの文数・素性数・ラベル数の統計

4 実験

基本句チャンキングと固有表現抽出の2つのNLPタスクについてCoNLL2000 [11]およびCoNLL2002 [10]のデータセットを使い、提案するNewton-CG法(以降、NCG)の有効性を検証した。

実験では1次のCRFを使用した。チャンキングの素性はCRF++¹に付属の素性テンプレートと同じ2値素性を、固有表現抽出は文献[1]と同じ2値素性を使用した。表1にタスク毎の訓練・開発・テスト文数および素性数・ラベル数を示す。また、式(1)のハイパーパラメータ σ は訓練セットの分割または開発セットを用いて選択した。

比較対象として、バッチ学習にはLBFSGを、オンライン学習には確率的勾配法(SGD)を実装した。SGDの更新式は $\theta_{k+1} = \theta_k - \eta_k g_k$ である。ただし、 $\eta_k > 0$ は学習率であり、勾配 g_k はランダムに選択した1事例 $(x, y) \in D$ を用いて次のように定義する。

$$g_k = -\Phi(x, y) + \mathbb{E}_{P_{\theta_k}}(\Phi(x, y)) + \frac{\theta_k}{n\sigma^2}$$

なお、 $n \equiv |D|$ は全訓練文数である。学習率には k の逆数で減衰する式 $\eta_k = \frac{\eta_0}{1+k/n}$ (以降、SGD(1/k)) [3]と、指数的に減衰する式 $\eta_k = \eta_0 \alpha^{k/n}$ (以降、SGD(α^k)) [12]の2種を用いた。ただし、初期学習率 η_0 は最初に500文で学習し500文で f を評価して選択した。また、文献[12]より $\alpha = 0.85$ を使用した。さらに、疎な更新を効率的に行うためスカラー u とベクトル v で $\theta = uv$ を実装した [9]。

実装はすべてJava™で行い、Java HotSpot™ 64-Bit Server VM (1.6.0_18)上で実行した。CPUはIntel® Core™2 Quad 2.4 GHz、メモリは8 GBである。

図1に、チャンキング・タスクにおける時間毎(横軸)の目的関数の最小値 f_{\min} との差(縦軸)を示す。なお、固有表現抽出タスクの傾向は同じであるため省略する。図中のNCGの c の値は3節で示した周辺確率を保存した文数を示す($c=ALL$ は全文の周辺確率を保存した結果)。また、LBFSGの m はヘシアン逆行列を近似するために差分ベクトルを保存した履歴の長さであり、SGDは $k = 200n$ までの結果を示す。

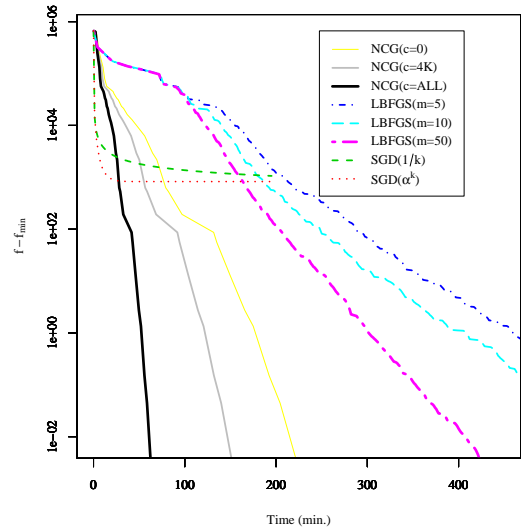


図 1: 目的関数の最小値との誤差(対数軸)と訓練時間

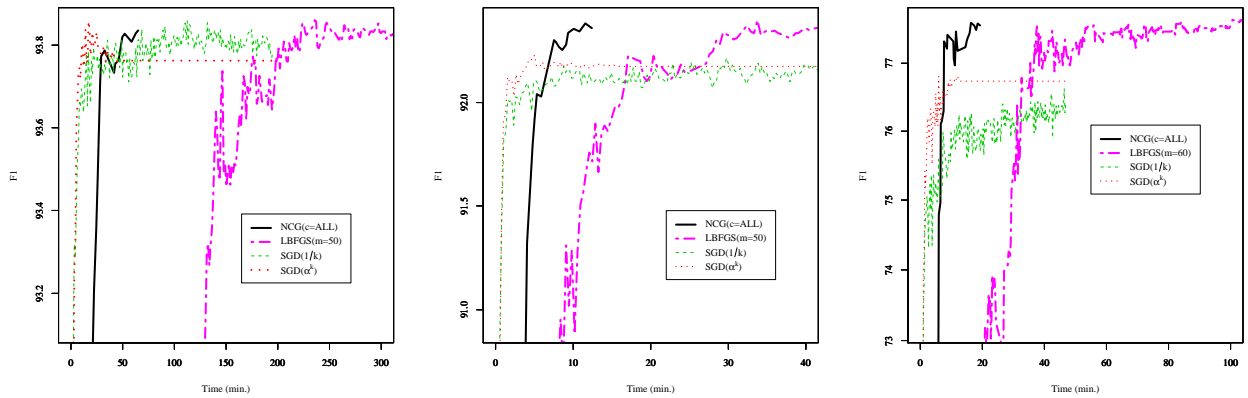
NCGもLBFSGも c および m を増やすことで収束が早くなっており、空間使用量と計算時間のトレードオフの関係があることがわかる。しかし、NCG($c=ALL$)はLBFSG($m = 50$)の6倍以上早く収束していることがわかる。なお、搭載メモリの制約からLBFSGの履歴は $m = 50$ 以上増やすことができなかった。この差は、LBFSGはNewton方向を近似するための情報量が固定であるのに対して、NCGは適応的にNewton方向の正確さを上げていくことで効率よく最適化できていることに拠ると考えられる。また、SGDは初期に大きく目的関数を減少させるが、後半は減少が緩やかになる(SGD(1/k))または停止(SGD(α^k))した。

次に、図2にテスト文でのF1値(縦軸)の時間経過を示す。なお、小規模データでの性能を調べるため、チャンキングタスクのデータを1/4にした結果を図2(b)に示している。また、バッチ学習のNCGとLBFSGは最も速いNCG($c=ALL$)およびLBFSG($m = 50$)(チャンキング)、LBFSG($m = 60$)(固有表現抽出)の結果のみを示す。

図2よりすべての場合において、NCGが安定して早く高いF1値を達成できていることがわかる。LBFSGはすべての場合において最高のF1値を達成するためにNCGの数倍の時間がかかっている。また、SGDは学習の早い段階で高い性能を示すが、必ずしも最高のF1値を達成できていない。なお、チャンキング全データ(図2(a))においてSGD(α^k)はNCGより早く最高F1値を達成しているが、学習の後半では低いF1値に収束してしまっている。その他の結果でもSGD(α^k)は早い段階で学習が停まっており、学習率の減衰が早すぎることをわかる。

以上の実験結果より、本実験で使用した規模の学習

¹<http://crfpp.sourceforge.net/>



(a) 基本句チャンキング

(b) 基本句チャンキング (サイズ 1/4)

(c) 固有表現抽出

図 2: テストデータでの F1 値と訓練時間

タスクでは提案法が最も安定して高速であることが示された。ただし、紙面の都合で省略するが、訓練データ約 4 万文の Penn Treebank の品詞タグ付けタスク [6] では提案法より SGD が最高性能を早く達成できた。しかし、多くの NLP のタスクでは大規模なコーパスを用意することは容易ではなく、提案法は中規模までの NLP タスクの学習に向いているといえる。

5 結論

本研究では、CRF のバッチ学習の高速化により、バッチ学習が有効となる学習データ規模を引き上げることに成功した。小・中規模タスクでは、更新率の調整や停止条件などに課題のあるオンライン学習手法よりも高速で安定した手法であることを実験により示した。今後は、スパース学習やオンライン学習への拡張などを検討したい。

参考文献

- [1] Yasemin Altun, Mark Johnson, and Thomas Hofmann. Investigating loss functions and optimization methods for discriminative learning of label sequences. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 145–152, 2003.
- [2] Minmin Chen, Yixin Chen, and Michael R. Brent. CRF-OPT: An efficient high-quality conditional random field solver. In *Proceedings of 23rd AAAI Conference on Artificial Intelligence*, pp. 1018–1023, 2008.
- [3] Michael Collins, Amir Globerson, Terry Koo, Xavier Carreras, and Peter L. Bartlett. Exponentiated gradient algorithms for conditional random fields and max-margin markov networks. *Journal of Machine Learning Research*, Vol. 9, pp. 1775–1822, 2008.
- [4] John Lafferty, Andrew McCallum, and Fernando Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the 18th International Conference on Machine Learning*, pp. 282–289, 2001.
- [5] Chih-Jen Lin, Ruby C. Weng, and S. Sathiyaraj. Trust region Newton method for large-scale logistic regression. *Journal of Machine Learning Research*, Vol. 9, pp. 627–650, 2008.
- [6] Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. Building a large annotated corpus of English: The Penn treebank. *Computational Linguistics*, Vol. 19, No. 2, pp. 313–330, 1993.
- [7] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer-Verlag, second edition, 2006.
- [8] Fei Sha and Fernando Pereira. Shallow parsing with conditional random fields. In *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology*, pp. 134–141, 2003.
- [9] Shai Shalev-Shwartz, Yoram Singer, and Nathan Srebro. Pegasos: Primal Estimated sub-GrAdient SOLver for svm. In *Proceedings of the 24th International Conference on Machine Learning*, pp. 807–814, 2007.
- [10] Erik F. Tjong Kim Sang. Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition. In *Proceedings of CoNLL-2002*, pp. 155–158, 2002.
- [11] Erik F. Tjong Kim Sang and Sabine Buchholz. Introduction to the CoNLL-2000 shared task: Chunking. In *Proceedings of CoNLL-2000*, pp. 127–132, 2000.
- [12] Yoshimasa Tsuruoka, Jun'ichi Tsujii, and Sophia Ananiadou. Stochastic gradient descent training for L1-regularized log-linear models with cumulative penalty. In *Proceedings of ACL-IJCNLP*, pp. 477–485, 2009.