

フランク・ウルフ法に基づく 1-ノルム正則化ソフトマージン最適化

[春の OR 学会]

© 九州大学/理研 AIP

三星 諒太郎

九州大学/理研 AIP

畑埜 晃平

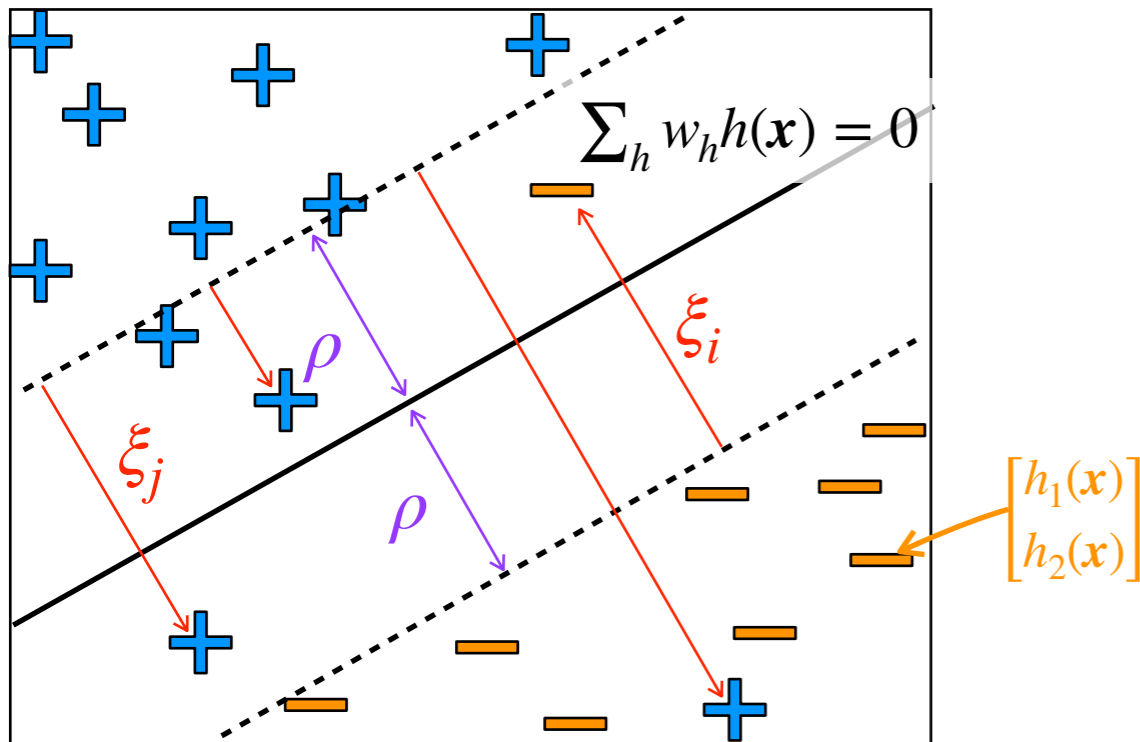
九州大学

瀧本 英二

「ソフト」マージン最適化

$H \subset [-1, +1]^X$ 中の関数の凸結合 $\text{sign} \circ \sum_{h \in H} w_h h$ で

サンプル $\{(x_i, y_i)\}_{i=1}^m \subset X \times \{\pm 1\}$ を分類.



「ソフト」マージン最適化問題

$$\begin{aligned} \max_{\rho, w, \xi} \quad & \rho - \frac{1}{\nu} \sum_{i=1}^m \xi_i \\ \text{s.t.} \quad & y_i \sum_{h \in H} w_h h(x_i) \geq \rho - \xi_i, \quad \forall i = 1, 2, \dots, m \\ & \sum_{h \in H} w_h = 1, \quad w \geq \mathbf{0}, \xi \geq \mathbf{0}. \end{aligned}$$

- ✓ 高々 $\nu \leq m$ 個の点を外れ値とみなし，残りの点でマージン ρ の最大化を目指す.
- ✓ 未知のデータに対しても高い汎化性能を保証 [Bertlett, '98].
- ✓ 線形計画問題だが，一般に H は非常に大きいので直接解くのは非現実的.
 - ✓ H 中の「有用な」関数だけを使うことはできないか？
- ➔ ブースティング（列生成法）によって解く.

ブースティングの Protokol

ブースティング ... 「**ブースター**」と「**弱学習者**」との間の繰り返しゲーム。

各ラウンド $t = 1, 2, \dots, T$ において：

① S 上の分布 $d_t \in \Delta_{m,\nu}$ を決める。

ブースター

弱学習者

② 関数 $h_t \in H$ を返す。

$$\Delta_{m,\nu} = \{d \in [0, 1/\nu]^m \mid \|d\|_1 = 1\}$$

直感的には

- **ブースター**は $\{h_1, h_2, \dots, h_{t-1}\}$ 上のどの仮説も正答率が小さくなる分布を選ぶ。
- **弱学習者**はサンプル上の分布 $d_t \in \Delta_{m,\nu}$ に対して正答率が一定より大きい h_t を返す。

$$\begin{aligned} & \max_{\rho, w, \xi} \rho - \frac{1}{\nu} \sum_{i=1}^m \xi_i \\ & \text{s.t. } y_i \sum_{h \in H} w_h h(x_i) \geq \rho - \xi_i, \quad \forall i = 1, 2, \dots, m \\ & \sum_{h \in H} w_h = 1, \quad w \geq \mathbf{0}, \xi \geq \mathbf{0}. \end{aligned}$$



$$\min_{d \in \Delta_{m,\nu}} \max_{h \in H} \sum_{i=1}^m d_i y_i h(x_i)$$

関連研究と主結果

	LPBoost [Demiriz+, 02]	ERLPBoost [Warmuth+, '08]	Cor. ERLPBoost [Shalev-Shwartz+, '10]	本研究
反復回数	$\Omega(m)$	$O\left(\frac{1}{\epsilon^2} \ln \frac{m}{\nu}\right)$	$O\left(\frac{1}{\epsilon^2} \ln \frac{m}{\nu}\right)$	$O\left(\frac{1}{\epsilon^2} \ln \frac{m}{\nu}\right)$
1 反復の計算	LP	CP	LP (Sorting)	\geq LP (Sorting)
実用面	非常に高速	遅い	ERLPB. より遅い	LPB. と同等

動機 : 理論保証があり, かつ実用上も高速なアルゴリズムが欲しい

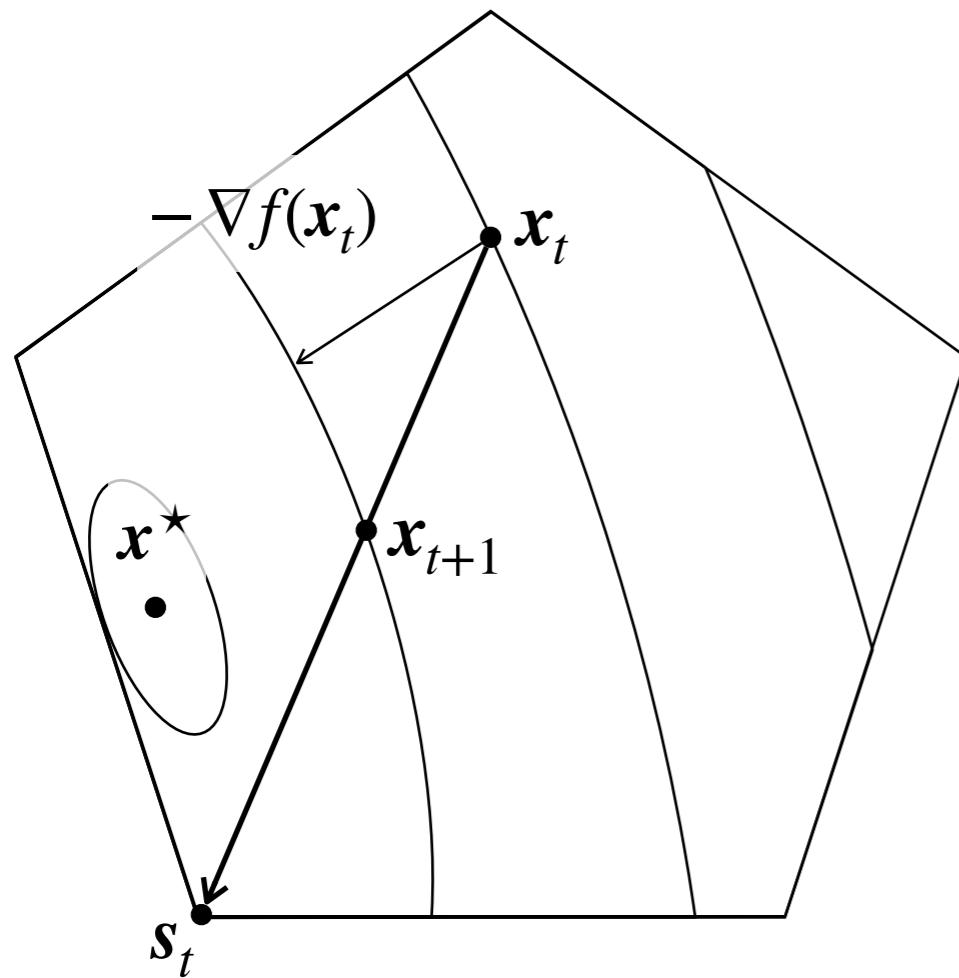
主結果

- 上記のアルゴリズム (後で詳述) は
すべて フランク・ウルフのアルゴリズムと見做せることを示した.
- 理論保証がある, **一般的な枠組**を提案.

フランク・ウルフのアルゴリズム [Marguerite+, '56]

閉凸集合上での凸関数の最小化を行う，一次の最適化アルゴリズム。

$$\min_{x \in S} f(x)$$



各ラウンド $t = 1, 2, \dots, T$ において：

1. $s_t \in \arg \max_{s \in S} s \cdot [-\nabla f(x_t)]$
2. $x_{t+1} = x_t + \lambda_t(s_t - x_t), \lambda_t \in [0, 1]$

λ_t の例

- $2/(t+1)$
- $f(x_{t+1})$ の x_t 周りの二次近似の最小解
- 線形探索

定理 [Jaggi, '13]

$$\forall x, y, \quad f(y) \leq f(x) + (y - x)^\top \nabla f(x) + \frac{\eta}{2} \|y - x\|^2$$

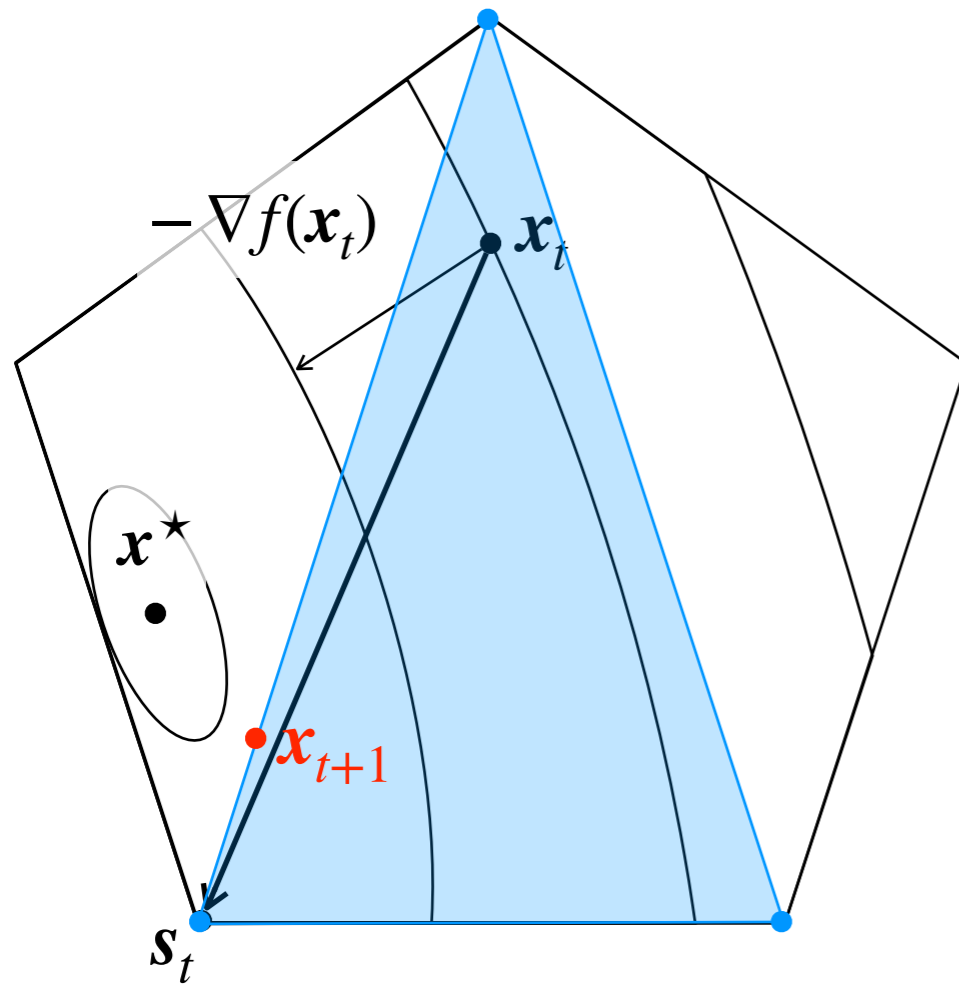
f が η -平滑であるとき，適切な $\{\lambda_t\}_t$ のもと

高々 $T = O(\eta/\epsilon)$ 回の反復で ϵ -近似解に収束。

フランク・ウルフのアルゴリズム [Marguerite+, '56]

閉凸集合上での凸関数の最小化を行う，一次の最適化アルゴリズム。

$$\min_{x \in S} f(x)$$



各ラウンド $t = 1, 2, \dots, T$ において：

1. $s_t \in \arg \max_{s \in S} s \cdot [-\nabla f(x_t)]$
2. $x_{t+1} = x_t + \lambda_t(s_t - x_t), \lambda_t \in [0, 1]$

その他の λ_t の例：FCFW

- 過去に求めた端点 $\{s_k\}_{k=1}^t \subset S$ の凸包上での f の最小解

定理 [Jaggi, '13]

$$\forall x, y, \quad f(y) \leq f(x) + (y - x)^\top \nabla f(x) + \frac{\eta}{2} \|y - x\|^2$$

f が η -平滑であるとき，適切な $\{\lambda_t\}_t$ のもと

高々 $T = O(\eta/\epsilon)$ 回の反復で ϵ -近似解に収束。

LPBoost [Demiriz+, '02] のアイデア

ソフトマージン最適化問題の双対問題を行生成法で解く

$$\begin{aligned} & \min_{\gamma, d \in \Delta_{m, \nu}} \gamma && \text{--- } (\spadesuit) \\ & \text{s.t. } \sum_{i=1}^m d_i y_i h(x_i) \leq \gamma, \quad \forall h \in H \end{aligned}$$

Lagrange
双対問題

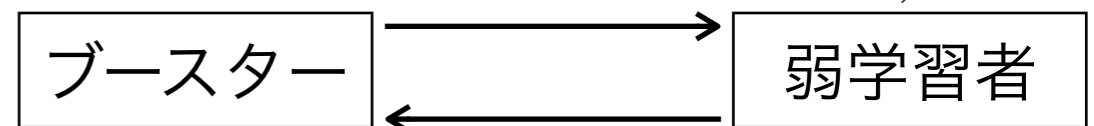
$$\begin{aligned} & \max_{\rho, w, \xi} \rho - \frac{1}{\nu} \sum_{i=1}^m \xi_i \\ & \text{s.t. } y_i \sum_{h \in H} w_h h(x_i) \geq \rho - \xi_i, \quad \forall i = 1, 2, \dots, m \\ & w \in \Delta_H, \quad \xi \geq 0. \end{aligned}$$

$$\Delta_H = \{w \in [0, 1]^H \mid \|w\|_1 = 1\}$$

- (\spadesuit) も H が大きいときは解くのが難しい。
- $H_0 = \emptyset$ から始め、 t ラウンド目において $H_{t-1} = \{h_1, h_2, \dots, h_{t-1}\}$ 上での (\spadesuit) の最適解 d_t を弱学習者に渡す。

各ラウンド $t = 1, 2, \dots, T$ において：

① H_{t-1} 上の (\spadesuit) の解 $d_t \in \Delta_{m, \nu}$.



② 関数 $h_t \in H$ を返す。

ERLPBoost [Warmuth+, '08] のアイデア

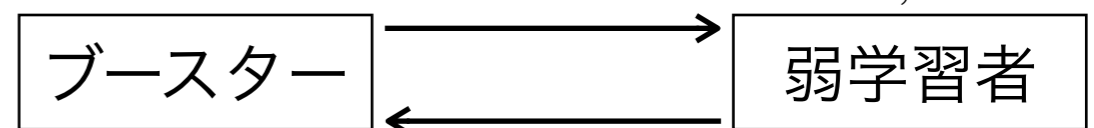
LPBoost + エントロピー正則化

$$\begin{aligned} \min_{\gamma, d \in \Delta_{m, \nu}} \quad & \gamma + \frac{1}{\eta} \sum_{i=1}^m d_i \ln(m d_i) \quad \text{--- } (\diamond) \\ \text{s.t.} \quad & \sum_{i=1}^m d_i y_i h(x_i) \leq \gamma, \quad \forall h \in H \end{aligned}$$

- $\eta > 0$ を十分大きく取れば, 正則化項は目的関数値にあまり影響を与えない
- $H_0 = \emptyset$ から始め, t ラウンド目において
 $H_{t-1} = \{h_1, h_2, \dots, h_{t-1}\}$ 上での (\diamond) の最適解 d_t を弱学習者に渡す.

各ラウンド $t = 1, 2, \dots, T$ において:

① H_{t-1} 上の (\diamond) の解 $d_t \in \Delta_{m, \nu}$.



② 関数 $h_t \in H$ を返す.

提案手法のアイデア

$A \in [-1, +1]^{m \times H}$ を, 第 (i, h) 成分が $A_{i,h} = y_i h(x_i)$ である行列とする.

$$\begin{aligned} \min_{\gamma, \mathbf{d} \in \Delta_{m,\nu}} \quad & \gamma + \frac{1}{\eta} \sum_{i=1}^m d_i \ln(m d_i) \\ \text{s.t.} \quad & \sum_{i=1}^m d_i y_i h(x_i) \leq \gamma, \quad \forall h \in H \end{aligned}$$

$$f(\mathbf{d}) \triangleq \begin{cases} 0, & \mathbf{d} \in \Delta_{m,\nu} \\ +\infty, & \mathbf{d} \notin \Delta_{m,\nu} \end{cases}$$

$$\tilde{f}(\mathbf{d}) \triangleq f(\mathbf{d}) + \frac{1}{\eta} \sum_{i=1}^m d_i \ln(m d_i)$$

① 制約を目的関数に

$$\min_{\mathbf{d}} \max_{h \in H} (\mathbf{d}^\top A)_h + \tilde{f}(\mathbf{d})$$

② Fenchel 双対問題

$$\max_{\mathbf{w} \in \Delta_H} -\tilde{f}^*(-A\mathbf{w})$$



$$\min_{\mathbf{w} \in \Delta_H} \tilde{f}^*(-A\mathbf{w})$$

③ フランク・ウルフ法で解く
最適解にのみ興味がある

Point!

• \tilde{f}^* はフランク・ウルフ法の仮定を満たす

f の Fenchel 共役関数 f^* : $f^*(\boldsymbol{\mu}) = \sup_{\boldsymbol{\theta}} \boldsymbol{\mu}^\top \boldsymbol{\theta} - f(\boldsymbol{\theta})$

提案手法のアイデア

フランク・ウルフ法

$$\min_{w \in \Delta_H} \tilde{f}^*(-Aw)$$

$$(Aw)_i = y_i \sum_{h \in H} w_h h(x_i)$$

$$\min_{x \in S} f(x)$$

1. 勾配ベクトルの計算

$$\begin{aligned} d_t &= \nabla \tilde{f}^*(-Aw_{t-1}) \\ &= \arg \max_d d^\top (-Aw_{t-1}) - \tilde{f}(d) \\ &= \arg \min_{d \in \Delta_{m,\nu}} d^\top Aw_{t-1} + \frac{1}{\eta} \sum_{i=1}^m d_i \ln(md_i) \end{aligned}$$

2. 端点の計算

$$\begin{aligned} h_t &\in \arg \min_{h \in H} d_t^\top (-Ae_h) \\ &= \arg \max_{h \in H} (d_t^\top A)_h \end{aligned}$$

arg max でなくともある値以上の h が得られるならば収束性を言える (後述)

3. 重みの更新

$$w_t = w_{t-1} + \lambda_t (e_{h_t} - w_{t-1})$$

通常の FW だけでなく、任意のアルゴリズムを組み込むことが可能！

各ラウンド $t = 1, 2, \dots, T$ において：

1. **Get** $\nabla f(x_{t-1})$
2. $s_t \in \arg \min_{s \in S} s \cdot \nabla f(x_{t-1})$
3. $x_t = x_{t-1} + \lambda_t (s_t - x_{t-1})$,
 $\lambda_t \in [0, 1]$

フランク・ウルフ法でブースティングを説明

ブースティング

$$\min_{\mathbf{w} \in \Delta_{m,\nu}} \tilde{f}^*(-A\mathbf{w})$$

各ラウンド $t = 1, 2, \dots, T$ において：

1. S 上の分布 $\mathbf{d}_t \in \Delta_{m,\nu}$ を決める。

ブースター \longleftrightarrow 弱学習者

2. 関数 $h_t \in H$ を返す。

例：ERLPBoost [Warmuth+, '08]

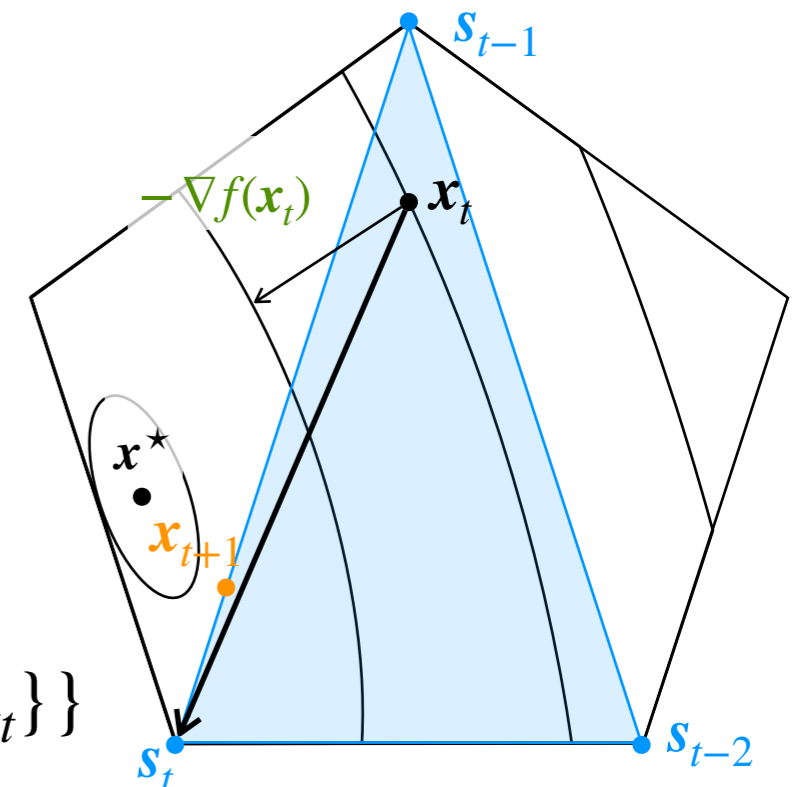
1. $\mathbf{d}_t = \arg \min_{\mathbf{d}} \max_{h \in \{h_1, h_2, \dots, h_{t-1}\}} (\mathbf{d}^\top A)_h + \tilde{f}(\mathbf{d})$
 $= \nabla \tilde{f}^*(\boldsymbol{\theta}_{t-1})$, 但し $\boldsymbol{\theta}_{t-1} = -A\mathbf{w}_{t-1}$
2. $\mathbf{e}_{h_t} \in \arg \min_{\mathbf{e}_h: h \in H} \mathbf{d}^\top (-A\mathbf{e}_h) = \arg \max_{\mathbf{e}_h: h \in H} \mathbf{d}^\top A\mathbf{e}_h$
3. $\mathbf{w}_t = \arg \min_{\mathbf{w} \in \Delta_{H_t}} \tilde{f}^*(-A\mathbf{w})$,
 但し $\Delta_{H_t} = \{\mathbf{w} \in \Delta_H \mid w_h > 0 \implies h \in \{h_1, h_2, \dots, h_t\}\}$

フランク・ウルフ法

$$\min_{\mathbf{x} \in S} f(\mathbf{x})$$

各ラウンド $t = 1, 2, \dots, T$ において：

1. Get $\nabla f(\mathbf{x}_{t-1})$
2. $\mathbf{s}_t \in \arg \max_{\mathbf{s} \in S} \mathbf{s} \cdot [-\nabla f(\mathbf{x}_{t-1})]$
3. $\mathbf{x}_t = \mathbf{x}_{t-1} + \lambda_t(\mathbf{s}_t - \mathbf{x}_{t-1})$, $\lambda_t \in [0, 1]$



フランク・ウルフ法でブースティングを説明

ブースティング

$$\min_{\mathbf{w} \in \Delta_{m,\nu}} f^*(-A\mathbf{w})$$

各ラウンド $t = 1, 2, \dots, T$ において：

1. S 上の分布 $\mathbf{d}_t \in \Delta_{m,\nu}$ を決める。

ブースター

弱学習者

2. 関数 $h_t \in H$ を返す。

フランク・ウルフ法

$$\min_{\mathbf{x} \in S} f(\mathbf{x})$$

各ラウンド $t = 1, 2, \dots, T$ において：

1. Get $\nabla f(\mathbf{x}_{t-1})$
2. $\mathbf{s}_t \in \arg \max_{\mathbf{s} \in S} \mathbf{s} \cdot [-\nabla f(\mathbf{x}_{t-1})]$
3. $\mathbf{x}_t = \mathbf{x}_{t-1} + \lambda_t(\mathbf{s}_t - \mathbf{x}_{t-1}), \lambda_t \in [0, 1]$

例：LPBoost [Demiriz+, '02]

$$1. \mathbf{d}_t \in \arg \min_{\mathbf{d}} \max_{h \in \{h_1, h_2, \dots, h_{t-1}\}} (\mathbf{d}^\top A)_h + f(\mathbf{d}) \\ = \nabla \tilde{f}^*(\boldsymbol{\theta}_{t-1}), \text{ 但し } \boldsymbol{\theta}_{t-1} = -A\mathbf{w}_{t-1}$$

$$2. \mathbf{e}_{h_t} \in \arg \min_{\mathbf{e}_h: h \in H} \mathbf{d}^\top (-A\mathbf{e}_h)$$

$$3. \mathbf{w}_t = \arg \min_{\mathbf{w} \in \Delta_{H_t}} f^*(-A\mathbf{w}), \\ \text{但し } \Delta_{H_t} = \{\mathbf{w} \in \Delta_H \mid w_h > 0 \implies h \in \{h_1, h_2, \dots, h_t\}\}$$

$$f(\mathbf{d}) = \begin{cases} 0 & \mathbf{d} \in \Delta_{m,\nu} \\ +\infty & \mathbf{d} \notin \Delta_{m,\nu} \end{cases}$$

**f^* は平滑性を持たないので、
フランク・ウルフの保証を使えない。**

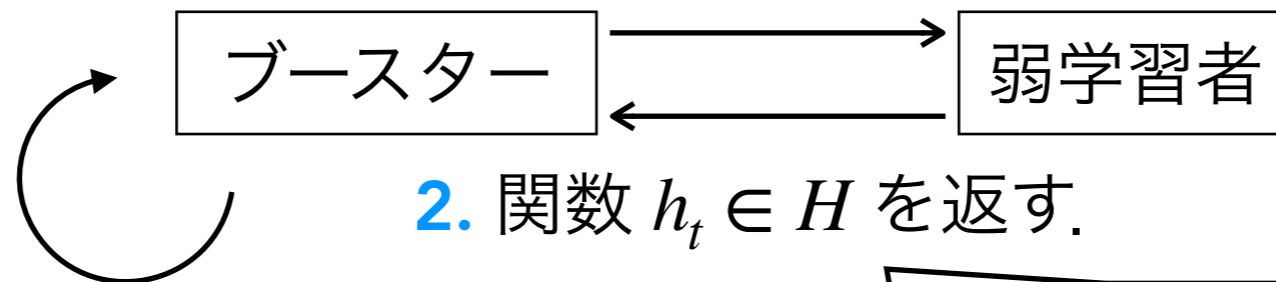
提案手法

フランク・ウルフ法の解析をよくみると次のことがわかる。

観察 各ラウンド t において

$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{2}{t+1}(\mathbf{s}_t - \mathbf{x}_{t-1})$ よりも関数値が小さくなる $\hat{\mathbf{x}}_t$ を使った場合,
 \mathbf{x}_t を使った場合と同じ収束率が保証される。

1. $\mathbf{d}_t = \nabla \tilde{f}^*(-A\mathbf{w}_{t-1}) \in \Delta_{m,\nu}$



3-1. $\mathbf{w}_t^{\mathcal{F}} = \mathbf{w}_{t-1} + \frac{2}{t+1}(\mathbf{s}_t - \mathbf{w}_{t-1})$

3-2. 任意に $\mathbf{w}_t^{\mathcal{A}}$ を決める。

3-3. $\mathbf{w}_t^{\mathcal{F}}, \mathbf{w}_t^{\mathcal{A}}$ のうち, $\tilde{f}^* \circ (-A)$ をより小さくする方を \mathbf{w}_t として採用。

$\sum_{i=1}^m d_i y_i h(\mathbf{x}_i) \geq g$ を満たす h を返す

定理 提案手法は $T = O\left(\frac{2}{\epsilon^2} \ln \frac{m}{\nu}\right)$ 反復後 $-f^*(-A\mathbf{w}_T) \geq g - \epsilon$ を達成。

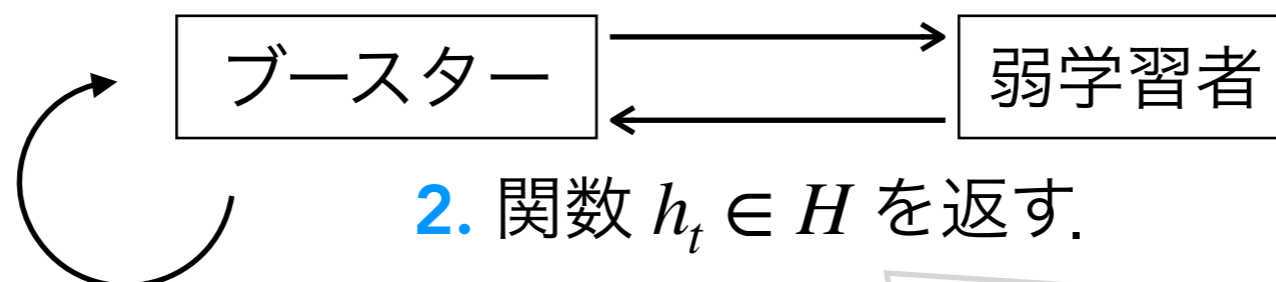
提案手法: 実用例

フランク・ウルフ法の解析をよくみると次のことがわかる.

観察 各ラウンド t において

$\mathbf{x}_t = \mathbf{x}_{t-1} + \frac{2}{t+1}(\mathbf{s}_t - \mathbf{x}_{t-1})$ よりも関数値が小さくなる $\hat{\mathbf{x}}_t$ を使った場合,
 \mathbf{x}_t を使った場合と同じ収束率が保証される.

$$1. \mathbf{d}_t = \nabla \tilde{f}^*(-A\mathbf{w}_{t-1}) \in \Delta_{m,\nu}$$



2. 関数 $h_t \in H$ を返す.

$$3-1. \mathbf{w}_t^{\mathcal{F}} = \mathbf{w}_{t-1} + \frac{2}{t+1}(\mathbf{s}_t - \mathbf{w}_{t-1})$$

3-2. 任意に $\mathbf{w}_t^{\mathcal{A}}$ を決める.

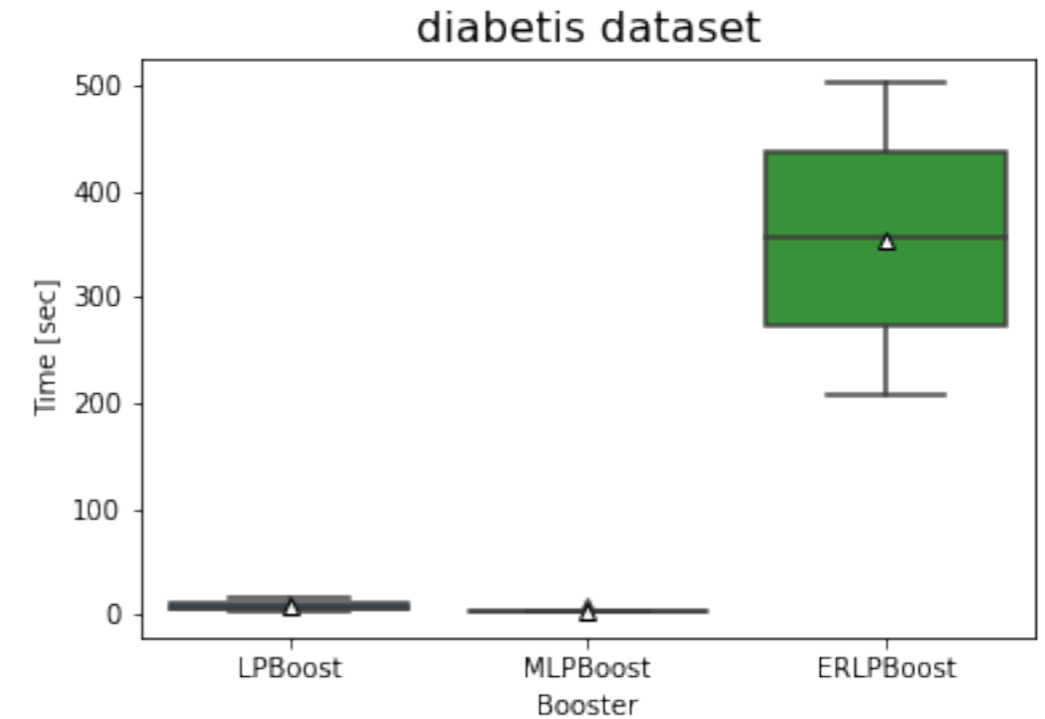
3-3. $\mathbf{w}_t^{\mathcal{F}}, \mathbf{w}_t^{\mathcal{A}}$ のうち, $\tilde{f}^* \circ (-A)$

- $\arg \max_{\mathbf{w} \in \Delta_{H_t}} -\tilde{f}^*(-A\mathbf{w}) \rightarrow \text{ERLPBoost}$
- LPBoost の双対解 $\rightarrow \text{MLPBoost (提案法)}$

実験

Setting.

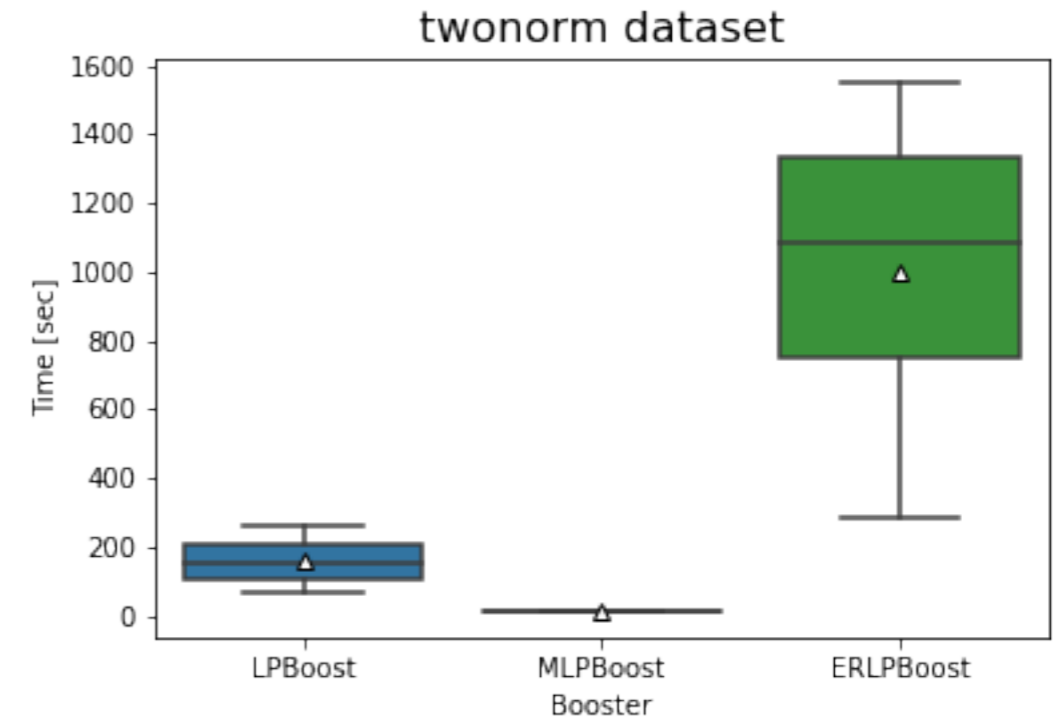
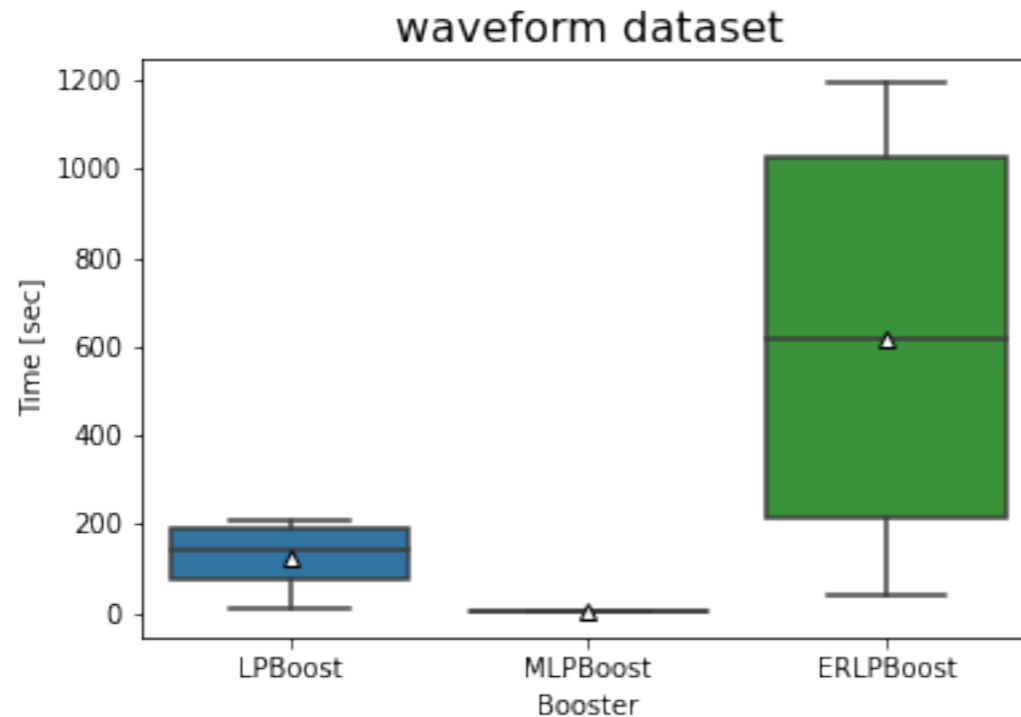
- 最適化ソルバ: **Gurobi 9.0.1**
- 弱学習器: 深さ 2 の決定木 (正答率基準)
- $\epsilon = 0.01$
- $\nu \in \{0.01, 0.02, 0.05, 0.1\}$ で計測



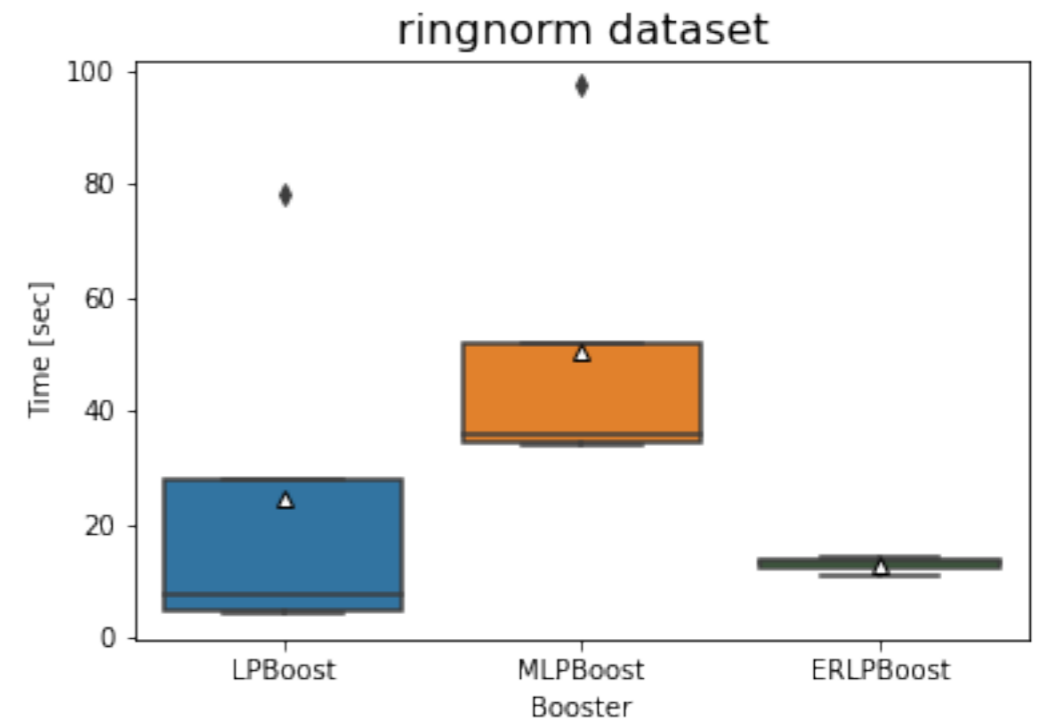
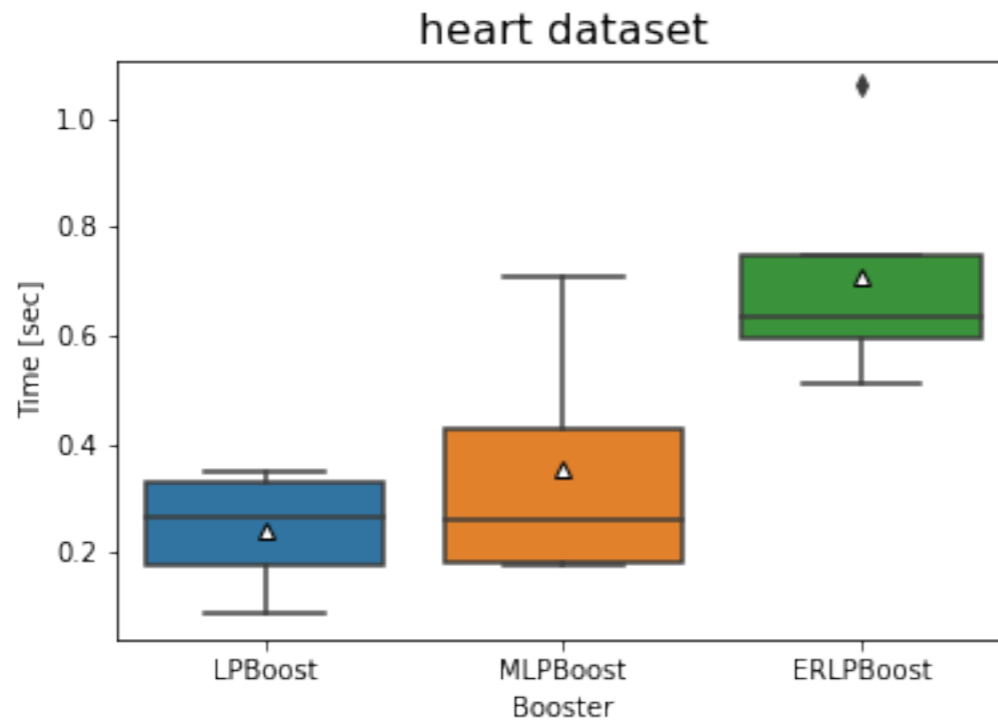
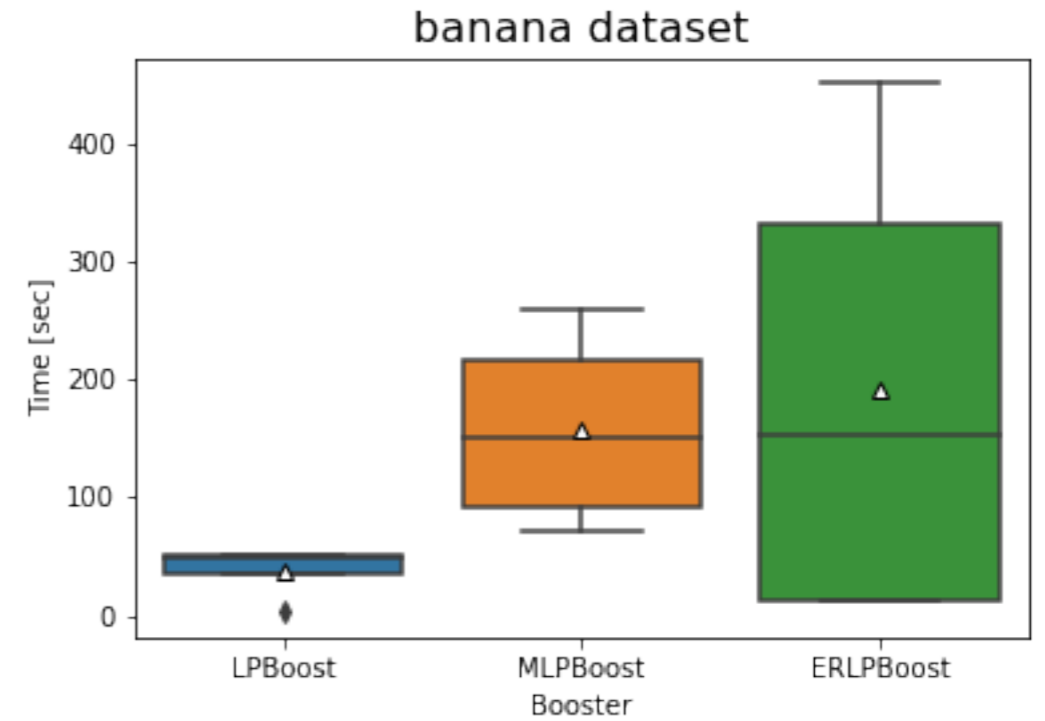
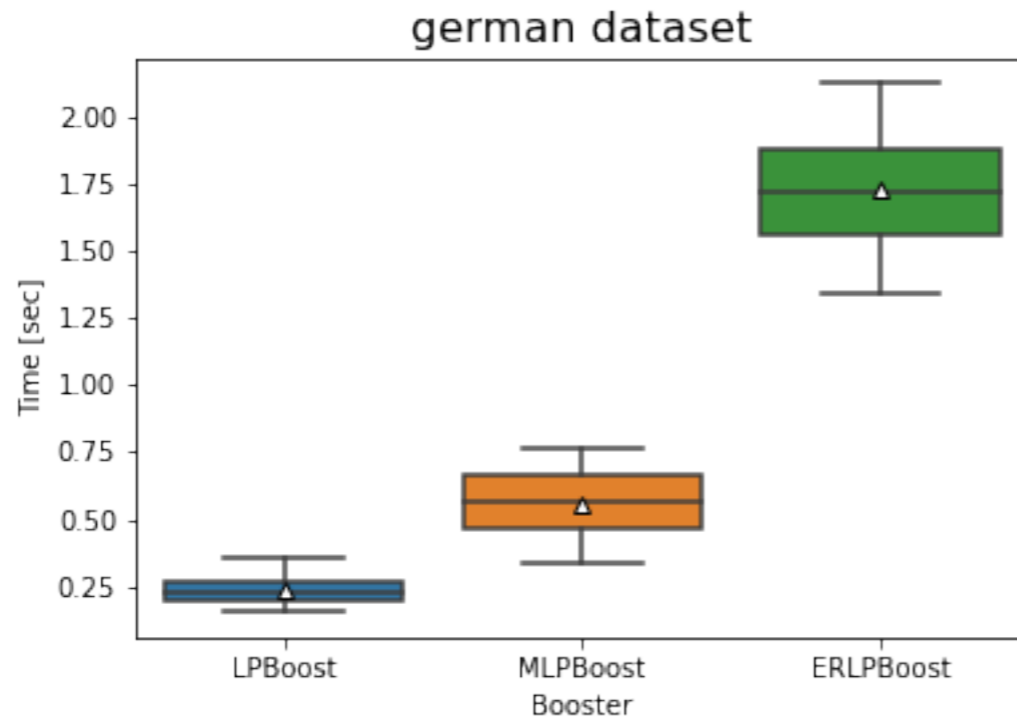
LPBoost

MLPBoost
(本研究)

ERLPBoost



実験



LPBoost
MLPBoost
(本研究)
ERLPBoost

まとめと今後の課題

- いくつかのブースティングアルゴリズムは
フランク・ウルフ法で説明できることを示した.
- 理論保証があり, かつ高速なブースティングアルゴリズムを提案.
 - **LPBoost** より (常に) 速くできるか?
- 強凸性を仮定すれば線形収束するか?