

PRML 第14章

宮澤 彬

総合研究大学院大学 博士前期

`miyazawa-a@nii.ac.jp`

September 25, 2015
(modified: October 4, 2015)

はじめに

- ▶ このスライドの Lua \LaTeX のソースコードは <https://github.com/pecorarista/documents> にあります.
- ▶ 教科書とは若干異なる表記をしている場合があります.
- ▶ 図は TikZ で描いたものを除き, Bishop (2006) のサポートページ <http://research.microsoft.com/en-us/um/people/cmbishop/prml/> または Murphy (2012) のサポートページ <http://www.cs.ubc.ca/~murphyk/MLbook/> から引用しています.
- ▶ 間違った記述を見つけた方は連絡ください.

今日やること

1. ベイズモデル平均化
2. コミッティ
3. ブースティング (AdaBoost)
 - ▶ 分類
 - ▶ ~~回帰~~
4. 木構造モデル (CART)
 - ▶ 回帰
 - ▶ 分類
5. ~~条件付き混合モデル~~

モデルの結合

複数のモデルの予測結果の平均を使って回帰を行うと、過学習が避けられて、より良い予測値が得られそう。また、分類の場合は多数決の結果を使うと同様のことが言えそう。

複数のモデルの出力を組み合わせて使う方法を**コミッティ**あるいは**合議** (committee) と呼ぶ。

モデルの結合

モデルの結合とベイズモデル平均化を混同しないように注意する.

まずはモデル結合の例として, 混合ガウス分布の密度推定を考える. モデルは, データがどの混合要素から生成されたかを示す潜在変数 z を使って, 同時確率 $p(x, z)$ で与えられる. 観測変数 x の密度は

$$p(x) = \sum_z p(x, z) \quad (14.3)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \quad (14.4)$$

となる. データ $X = (x_1, \dots, x_N)$ に関する周辺分布は

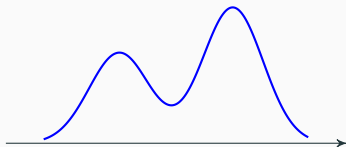
$$p(X) = \prod_{n=1}^N p(x_n) = \prod_{n=1}^N \left(\sum_{z_n} p(x_n, z_n) \right) \quad (14.5)$$

であり, 各 x_n に対して z_n が存在していることがわかる.

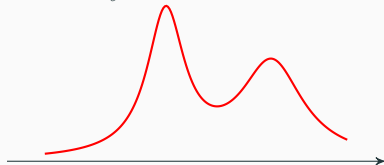
ベイズモデル平均化

次にベイズモデル平均化の例を示す。いくつかの異なるモデルがあって、 $h = 1, \dots, H$ で番号付けられているとする。例えば 1 つは混合ガウス分布で、もう 1 つは混合コーシー分布であるとする。

$$\sum_k \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)$$



$$\sum_{\ell} \rho_{\ell} t_1(x | \mu_{\ell}, \Sigma_{\ell})$$



ベイズモデル平均化

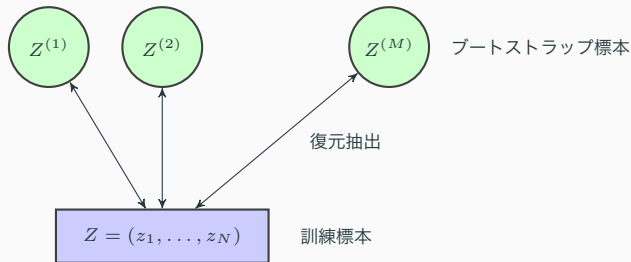
モデルに関する事前確率が $p(h)$ で与えられているならば、データに関する周辺分布は次式で求められる。

$$p(X) = \sum_{h=1}^H p(X|h) p(h) \quad (14.6)$$

モデル結合との違いは、データを生成しているのはモデル $h = 1, \dots, H$ のどれか 1 つだけということである。

コミッティ

コミッティを使う手法の例として、**バギング** (bagging; **b**ootstrap **a**ggregation) を紹介する。まず訓練用データ $Z = (z_1, \dots, z_N)$ から、ブートストラップ標本 $Z^{(1)}, \dots, Z^{(M)}$ を作る。



コミッティ

バギングにおけるコミッティの予測値 $y_{\text{COM}}(x)$ は、各ブートストラップ集合 $Z^{(1)}, \dots, Z^{(M)}$ で学習したモデルの予測 $y_1(x), \dots, y_M(x)$ を使って

$$y_{\text{COM}}(x) = \frac{1}{M} \sum_{m=1}^M y_m(x) \quad (14.7)$$

となる.

これで実際に良い結果が得られるのか、誤差を評価して確認する.

バギングの誤差評価

予測したい回帰関数を $h(x)$ とし、各モデルの予測値 $y_m(x)$ が、 $h(x)$ と加法的な誤差 $\varepsilon_m(x)$ を使って

$$y_m(x) = h(x) + \varepsilon_m(x) \quad (14.8)$$

と表せると仮定する。このとき、二乗誤差は

$$\mathbb{E} \left[(y_m(x) - h(x))^2 \right] = \mathbb{E} \left[(\varepsilon_m(x))^2 \right] \quad (14.9)$$

となるので、平均二乗誤差の期待値 E_{AV} は

$$\begin{aligned} E_{AV} &:= \mathbb{E} \left[\frac{1}{M} \sum_{m=1}^M (y_m(x) - h(x))^2 \right] \\ &= \frac{1}{M} \sum_{m=1}^M \mathbb{E} \left[(\varepsilon_m(x))^2 \right] \end{aligned} \quad (14.10)$$

で求められる。

バギングの誤差評価

一方、コミッティによる出力について、誤差の期待値 E_{COM} は

$$\begin{aligned} E_{\text{COM}} &:= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M y_m - h(x) \right)^2 \right] \\ &= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \varepsilon_m(x) \right)^2 \right] \end{aligned} \quad (14.11)$$

と求められる. ベクトル $(1 \cdots 1)' \in \mathbb{R}^M$ と $(\varepsilon_1(x) \cdots \varepsilon_M(x))' \in \mathbb{R}^M$ に対して Cauchy-Schwarz の不等式を用いると

$$\begin{aligned} \left(\sum_{m=1}^M \varepsilon_m(x) \right)^2 &\leq M \sum_{m=1}^M (\varepsilon_m(x))^2 \\ E_{\text{COM}} &\leq E_{\text{AV}} \end{aligned}$$

が成り立つ. つまりコミッティを使う場合の誤差の期待値は, 構成要素の誤差の期待値の和を超えない.

バギングの誤差評価

特に誤差の平均が 0 で無相関である, すなわち

$$\mathbb{E}[\varepsilon_m(x)] = 0 \quad (14.12)$$

$$\text{cov}(\varepsilon_m(x), \varepsilon_\ell(x)) = \mathbb{E}[\varepsilon_m(x) \varepsilon_\ell(x)] = 0, m \neq \ell \quad (14.13)$$

が成り立つならば¹

$$\begin{aligned} E_{\text{COM}} &= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \varepsilon_m(x) \right)^2 \right] \\ &= \frac{1}{M^2} \mathbb{E} \left[\sum_{m=1}^M (\varepsilon_m(x))^2 \right] \\ &= \frac{1}{M} E_{\text{AV}} \end{aligned} \quad (14.14)$$

となり, 誤差の期待値が大幅に低減される.

¹ 同じようなモデルを, 同じような訓練データで学習させているのだから, こんなことは期待できない.

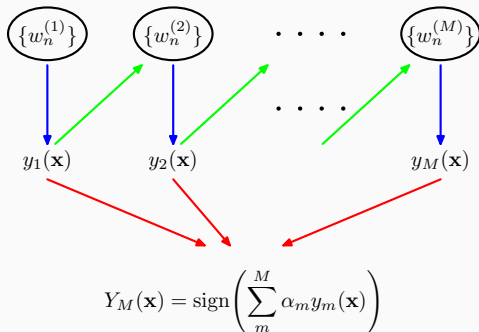
ブースティング

コミッティを使う手法の別な例として、**ブースティング** (boosting) と呼ばれる手法を紹介する。ブースティングは分類のためのアルゴリズムとして設計されたが、回帰にも拡張できる。とりあえず分類について考える。

コミッティを構成する個々の分類器は**ベース学習器** (base learner) または**弱学習器** (weak learner) と呼ばれる。ブースティングは、弱学習器を（並列的にではなく）逐次的に、重み付きデータを使って学習を進める方法である。重みの大きさは、それ以前に誤って分類されたデータに対して大きくなるように定める。

ブースティングの中で代表的な手法である **AdaBoost** (adaptive boosting) を紹介する。

AdaBoost



for $m = 1, \dots, N$

Initialize the data weighting coefficients

$$w_n^{(1)} := \frac{1}{N}$$

AdaBoost

for $m = 1, \dots, M$

Fit a classifier $y_m(x)$ to the training data by minimizing

$$J_m := \sum_{n=1}^N w_n^{(m)} \mathbb{1}_{\{x \mid y_m(x) \neq t_n\}}(x_n). \quad (14.15)$$

Evaluate

$$\alpha_m := \log \frac{1 - \varepsilon_m}{\varepsilon_m} \quad (14.17)$$

where

$$\varepsilon_m := \frac{\sum_{n=1}^N w_n^{(m)} \mathbb{1}_{\{x \mid y_m(x) \neq t_n\}}(x_n)}{\sum_{n=1}^N w_n^{(m)}}. \quad (14.15)$$

Update the data weighting coefficients

$$w_n^{(m+1)} := w_n^{(m)} \exp \left(\alpha_m \mathbb{1}_{\{x \mid y_m(x_n) \neq t_n\}}(x_n) \right). \quad (14.18)$$

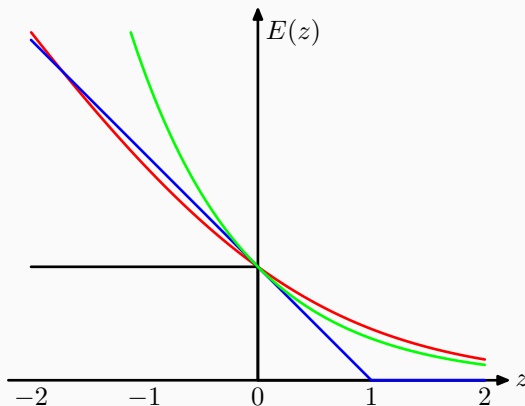
Make predictions using the final model, which is given by

$$Y_M(x) := \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(x) \right). \quad (14.19)$$

このアルゴリズム中に現れる式は、指数誤差の逐次的な最小化から導かれる。

指数誤差の最小化

目標値を $t \in \{-1, 1\}$, 予測値を $y \in \mathbb{R}$ とし, これらの積 ty を z と置く. このとき指数誤差は $E(z) := \exp(-z)$ という形で与えられる. $E(z)$ は t と y が同符号 ($z > 0$) なら 0 に近い正の値を, 異符号 ($z < 0$) ならば正の大きな値をとる.



指数誤差関数 交差エントロピー誤差 ヒンジ形誤差関数 誤分類誤差

指数誤差の最小化

以下の誤差 E を最小化したい.

$$E := \sum_{n=1}^N \exp \left(-t_n \sum_{\ell=1}^m \frac{1}{2} \alpha_{\ell} y_{\ell}(x_n) \right)$$

逐次的に学習していくので $\alpha_1, \dots, \alpha_{m-1}$ と $y_1(x), \dots, y_{m-1}(x)$ は所与として, α_m と $y_m(x)$ に関して E を最小化する. E を

$$\begin{aligned} E &= \sum_{n=1}^N \exp \left(-t_n \sum_{\ell=1}^{m-1} \frac{1}{2} \alpha_{\ell} y_{\ell}(x_n) - \frac{1}{2} t_n \alpha_m y_m(x_n) \right) \\ &= \sum_{n=1}^N w_n^{(m)} \exp \left(-\frac{1}{2} t_n \alpha_m y_m(x_n) \right) \end{aligned} \quad (14.22)$$

と書き換える. ここで以下のように置いた.

$$w_n^{(m)} := \exp \left(-t_n \sum_{\ell=1}^{m-1} \frac{1}{2} \alpha_{\ell} y_{\ell}(x_n) \right) \quad (14.22')$$

指数誤差の最小化

$y_m : x \mapsto \{-1, 1\}$ によって正しく分類されたデータ点の添字集合を \mathcal{T}_m , 誤って分類されたデータの添字集合を \mathcal{M}_m とする. このとき

$$\begin{aligned} E &= e^{-\alpha_m/2} \sum_{n \in \mathcal{T}_m} w_n^{(m)} + e^{\alpha_m/2} \sum_{n \in \mathcal{M}_m} w_n^{(m)} \\ &= \left(e^{\alpha_m/2} - e^{-\alpha_m/2} \right) \sum_{n=1}^N w_n^{(m)} \mathbb{1}_{\{x | y_m(x_n) \neq t_n\}}(x_n) \\ &\quad + e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} \end{aligned} \tag{14.23}$$

が成り立つ. これを $y_m(x)$ について最小化することは (14.15) を最小化することと同値である.

指数誤差の最小化

また α_m について微分して 0 と置くと

$$\left(e^{\alpha_m/2} + e^{-\alpha_m/2}\right) \sum_{n=1}^N w_n^{(m)} \mathbb{1}_{\{x|y_m(x_n) \neq t_n\}}(x_n) - e^{-\alpha_m/2} \sum_{n=1}^N w_n^{(m)} = 0$$

$$e^{\alpha_m} + 1 = \frac{\sum_{n=1}^N w_n^{(m)}}{\sum_{n=1}^N w_n^{(m)} \mathbb{1}_{\{x|y_m(x_n) \neq t_n\}}(x_n)}$$

となる. (14.15) と同じ置き換えを行うことにより

$$\alpha_m = \log \frac{1 - \varepsilon_m}{\varepsilon_m}$$

を得る. これはアルゴリズム中の (14.17) に等しい.

指数誤差の最小化

重みは (14.22') に従い次のように更新する.

$$w_n^{(m+1)} = w_n^{(m)} \exp \left(-\frac{1}{2} t_n \alpha_m y_m(x_n) \right) \quad (14.24)$$

ここで

$$t_n y_m(x_n) = 1 - 2\mathbb{1}_{\{x|y_m(x) \neq t_n\}}(x_n) \quad (14.25)$$

が成り立つことを使うと

$$\begin{aligned} w_n^{(m+1)} &= w_n^{(x)} \exp \left(-\frac{1}{2} \alpha_m (1 - 2\mathbb{1}_{\{x|y_m(x) \neq t_n\}}(x_n)) \right) \\ &= w_n^{(x)} \exp \left(-\frac{\alpha_m}{2} \right) \exp \left(\mathbb{1}_{\{x|y_m(x) \neq t_n\}}(x_n) \right) \end{aligned} \quad (14.26)$$

と書き換えられる. $\exp(-\alpha_m/2)$ はすべてのデータに共通なので無視すると (14.18) が得られる.

ブースティングのための誤差関数

以下の指数誤差を考える.

$$\mathbb{E} [\exp (-ty(x))] = \sum_{t \in \{-1,1\}} \int \exp (-ty(x)) p(t|x) p(x) dx \quad (14.27)$$

変分法を使って y について上式を最小化する.

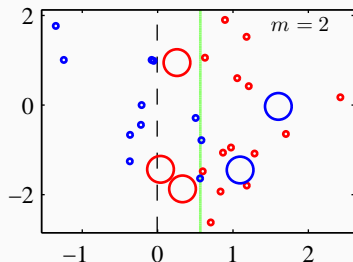
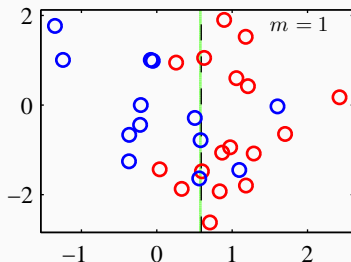
$\varphi_t(y) := \exp (-ty(x)) p(t|x) p(x)$ とする. 停留点は以下のように求められる.

$$\begin{aligned} \sum_{t \in \{-1,1\}} D_y \varphi_t(y) &= 0 \\ \exp(y(x)) p(t=-1|x) - \exp(-y(x)) p(t=1|x) &= 0 \\ y(x) &= \frac{1}{2} \log \frac{p(t=1|x)}{p(t=-1|x)} \end{aligned} \quad (14.28)$$

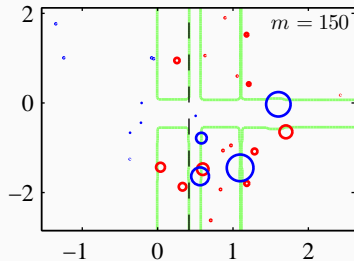
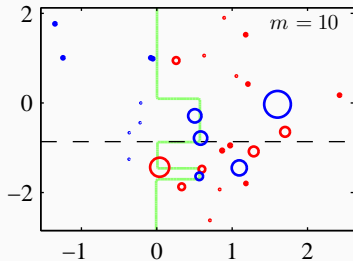
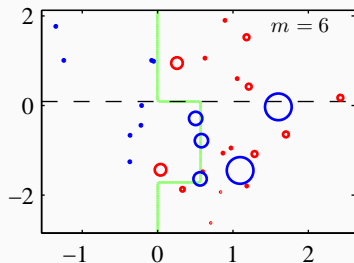
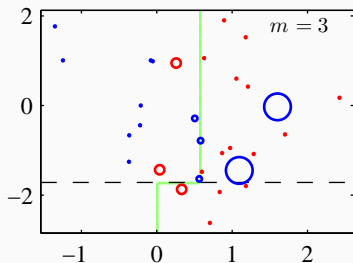
つまり AdaBoost は逐次的に $p(t=1|x)$ と $p(t=-1|x)$ の比の対数の近似を求めている. これが最終的に (14.19) で符号関数を使って分類器を作った理由となっている.

AdaBoost の振る舞い

AdaBoost による分類の例を示す。図中の円は各データを表し、その半径は割り振られた重みを表している。



AdaBoost の振る舞い



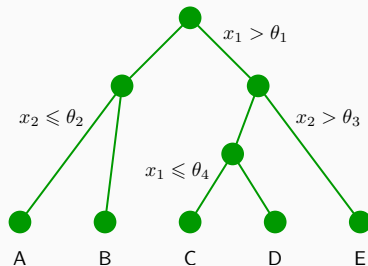
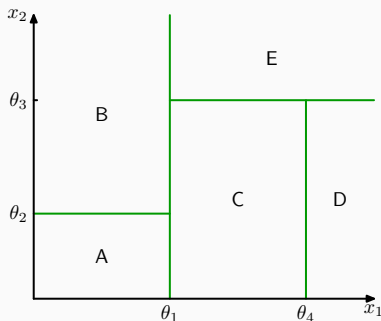
AdaBoost の拡張

AdaBoost は Freund and Schapire (1996) で導入されたが、指数誤差の逐次的最小化としての解釈は Friedman et al. (2000) で与えられた。これ以降、誤差関数の変更により様々な拡張が行われるようになった。

木構造モデル

今までは複数の出力を合わせる方法を見てきた。今度は入力を複数のモデルの 1 つに振り分ける方法である木構造モデルを扱う。

木構造モデルでは、入力空間を矩形 (cuboid) 領域に分けて、それぞれの領域に対応するモデルに予測をさせる。

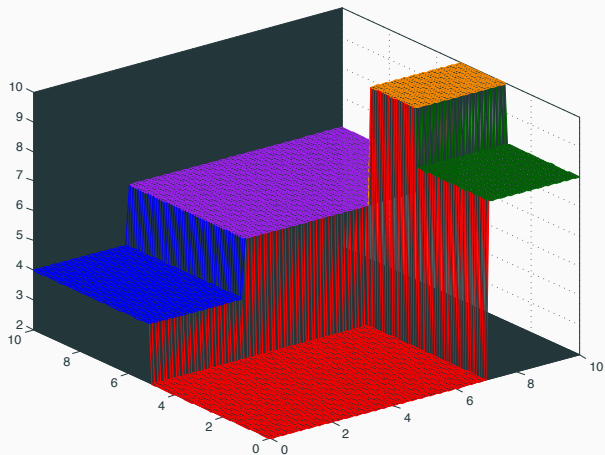


以下では木構造モデルの 1 つである **CART** (classification and regression tree) を扱う。CART は名前の通り分類にも回帰にも使うことができるが、まずは回帰に絞って説明する。

入力と目標値の組 $\mathcal{D} = (\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N) \in \mathbb{R}^D \times \mathbb{R}$ が与えられているとする。入力に対する分割が $\mathcal{R}_1, \dots, \mathcal{R}_M$ となっているとき、予測値を次のように定義する。

$$f(\mathbf{x}) := \sum_{m=1}^M c_m \mathbb{1}_{\mathcal{R}_m}(\mathbf{x})$$

CART



分割が与えられているときの最適な予測値

各領域で二乗和誤差を最小化する予測値を求めよう.

$\mathcal{I}_m := \{n; \mathbf{x}_n \in \mathcal{R}_m\}$ と表すと誤差は以下のようになる.

$$\begin{aligned}\sum_{n=1}^N (t_n - f(\mathbf{x}_n))^2 &= \sum_{n=1}^N \left(t_n - \sum_{m=1}^M c_m \mathbb{1}_{\mathcal{R}_m}(\mathbf{x}_n) \right)^2 \\ &= \sum_{m=1}^M \sum_{n \in \mathcal{I}_m} (c_m - t_n)^2\end{aligned}$$

各 c_m について微分して 0 と置くと, 最適な予測値は

$$\bar{c}_m = \frac{1}{|\mathcal{I}_m|} \sum_{n \in \mathcal{I}_m} t_n$$

となることが分かる. ただし二乗和誤差を最小化するような分割を決めようとするすると計算量が大きくなる. そこで**貪欲法**を使って木の構造を定める.

分割基準

各回の領域の分割基準を考える. 変数を表す添字を $j \in \{1, \dots, D\}$, 閾値を θ として, 2つの領域 $\mathcal{R}_1(j, \theta)$ と $\mathcal{R}_2(j, \theta)$ を次で定める.

$$\mathcal{R}_1(j, \theta) = \{\mathbf{x}; x_j \leq \theta\}, \quad \mathcal{R}_2(j, \theta) = \{\mathbf{x}; x_j > \theta\}$$

$\mathcal{I}_i(j, \theta) := \{n; \mathbf{x}_n \in \mathcal{R}_i(j, \theta)\}$ とする. 分割の基準となる (j, θ) は次を解いて得られる.

$$\min_{j \in \{1, \dots, D\}} \min_{\theta} \left\{ \min_{c_1} \sum_{n \in \mathcal{I}_1(j, \theta)} (c_1 - t_n)^2 + \min_{c_2} \sum_{n \in \mathcal{I}_2(j, \theta)} (c_2 - t_n)^2 \right\}.$$

内側の和を最小化する c_1 と c_2 は, それぞれ

$$\bar{c}_1 = \frac{1}{|\mathcal{I}_1(j, \theta)|} \sum_{n \in \mathcal{I}_1(j, \theta)} t_n, \quad \bar{c}_2 = \frac{1}{|\mathcal{I}_2(j, \theta)|} \sum_{n \in \mathcal{I}_2(j, \theta)} t_n$$

となる. あとは各 j について最適な θ を計算し, そのあとで最適な j を求めればよい.

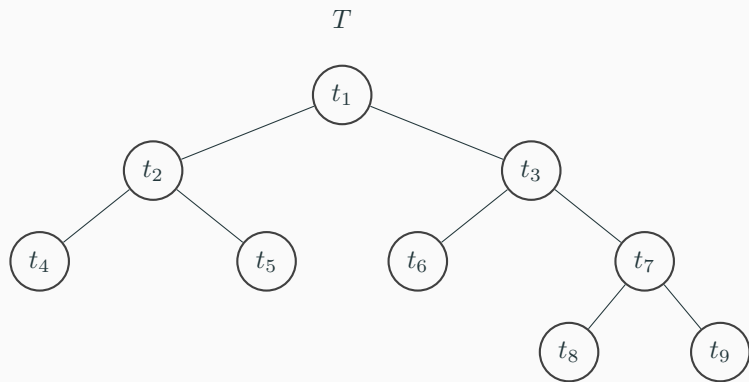
停止基準

次に停止基準について考えなければならない.

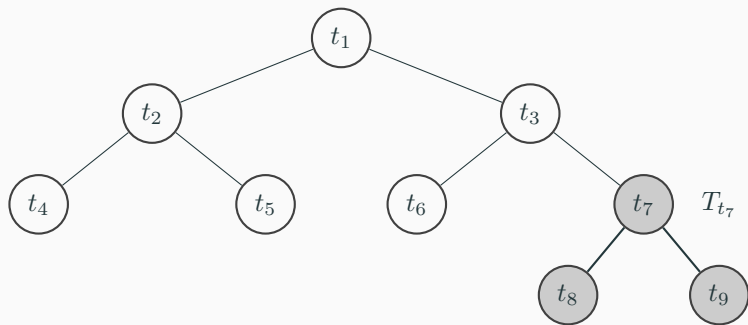
単純に考えると誤差の減少幅が小さくなったときに止める方法が良さそうであるが, 分割を続けていると大きな誤差の減少が見られることが経験的に知られている.

そこで, 先に大きな木を作ってからその木の枝を刈って適当な木を得ることにする.

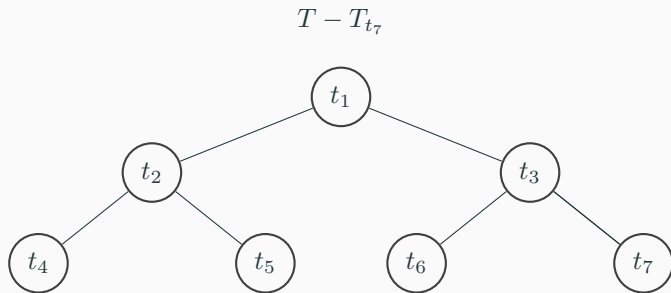
枝刈り



枝刈り



枝刈り



t_8 に対応する領域と t_9 に対応する領域を併合するので t_7 は残る.

枝刈りの基準

切り落とす位置は、次の **error-complexity measure**

$$R_{\alpha}(T) = R(T) + \alpha |\tilde{T}|$$

が一番大きく減少するような節点 t を選ぶ. ただし \tilde{T} は木 T の葉全体の集合, $R(T)$ は誤差

$$R(T) := \sum_{m=1}^{|\tilde{T}|} \sum_{n \in \mathcal{I}_m} (\bar{c}_m - t_n)^2$$

である. $\alpha \geq 0$ は枝の数に対する罰則を調整するパラメータである. この基準では, 刈っても (通常大きくなってしまう) 誤差があまり大きくならず, 葉の数が減るような枝刈りが優先して行われる.

枝刈りの基準

α を動かすと次のようなことがわかる.

- ▶ $\alpha = 0$ のとき
各データに 1 つの葉を対応させるような木が最適になる.
- ▶ α が大きいとき
根 t_1 だけから成る木が最適になる.

ちょうどいい木はこの間にある.

枝刈り

部分木 T_t と節点 t から下を刈ってできる木 $(\{t\}, \emptyset)$ の error-complexity measure はそれぞれ

$$R_\alpha(T_t) = R(T_t) + \alpha |\tilde{T}_t|,$$
$$R_\alpha((\{t\}, \emptyset)) = R((\{t\}, \emptyset)) + \alpha$$

である. $R_\alpha(T_t) > R_\alpha((\{t\}, \emptyset))$, すなわち

$$\alpha < \frac{R(T_t) - R((\{t\}, \emptyset))}{|\tilde{T}_t| - 1}$$

が成り立つならば刈り取ったほうがよい. この右辺を刈り取る効果の指標とし

$$g(t; T) := \frac{R(T_t) - R((\{t\}, \emptyset))}{|\tilde{T}_t| - 1}$$

と定める.

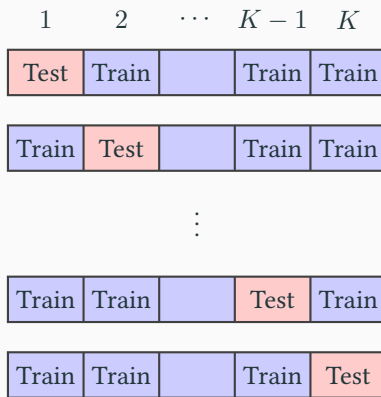
枝刈り

次の**最弱リンク枝刈り** (weakest link pruning) と呼ばれるアルゴリズムにより $T^0 \succ \dots \succ T^J = (\{t_1\}, \emptyset)$ と狭義単調増加列 $\alpha_0 < \dots < \alpha_J$ を得る.

```
1:  $i \leftarrow 0$   
2: while  $|\widetilde{T^i}| > 1$   
3:    $\alpha_i \leftarrow \min_{t \in V(T^i) \setminus \widetilde{T^i}} g(t; T^i)$   
4:    $\mathcal{T}^i \leftarrow \arg \min_{t \in V(T^i) \setminus \widetilde{T^i}} g(t; T^i)$   
5:   for  $t \in \mathcal{T}^i$   
6:      $T^i \leftarrow T^i - T_t^i$   
7:    $i \leftarrow i + 1$ 
```

交差確認

最終的にやらなければならないことは、得られた列 $\langle T^i \rangle_{i=0}^J$ の中から適当な T^i 選ぶことである。そのための方法としては交差確認 (cross validation) が使われる。



手順は以下のようになる.

- 1: データ \mathcal{D} を使って, 木の列 $\langle T^i \rangle_{i=0}^J$ とパラメータの列 $\langle \alpha_i \rangle_{i=0}^J$ を作る.
- 2: データ \mathcal{D} を $\mathcal{D} = \coprod_{k=1}^K \mathcal{D}_k$ と分割する. 極力, 各 $|\mathcal{D}_k|$ が同じ大きさになるようにする.
- 3: 各 $\mathcal{D}^{(k)} := \mathcal{D} \setminus \mathcal{D}_k$ を使い, 木の列 $\langle T^{(k)i} \rangle_{i=0}^{i_k}$ とパラメータの列 $\langle \alpha_i^{(k)} \rangle_{i=0}^{i_k}$ を作る.

- 4: 各 $\alpha'_i := \sqrt{\alpha_i \alpha_{i+1}}$ に対して $R_{\alpha'_i}(T)$ を最小化するような $T^{(1)}(\alpha'_i), \dots, T^{(K)}(\alpha'_i)$ とそれらに対応する予測 $y_i^{(1)}, \dots, y_i^{(K)}$ を求める. ここで $\alpha \in [\alpha_i^{(k)}, \alpha_{i+1}^{(k)}]$ ならば $T^{(k)}(\alpha) = T^{(k)i}$ となることに注意する.
- 5: 今までの結果から以下を求めることができる.

$$R^{\text{CV}}(T^i) = \frac{1}{N} \sum_{k=1}^K \sum_{n \in \{n' ; (\mathbf{x}_n, t_n) \in \mathcal{D}_k\}} \left(t_n - y_i^{(k)}(\mathbf{x}_n) \right)^2$$

- 6: 誤差が一番小さい木を $T^{**} := \arg \min_{T^i} R(T^i)$ とする.

7: 標準誤差 (standard error) SE を求める.

$$\begin{aligned} \text{SE} (R^{\text{CV}} (T^i)) &:= s (T^i) / \sqrt{N}, \\ s (T^i) &:= \sqrt{\frac{1}{N} \sum_{n=1}^N \left(\left(t_n - y_i^{(\kappa(n))} (x_n) \right)^2 - R^{\text{CV}} (T^i) \right)^2}, \\ \kappa (n) &= \sum_{k=1}^K k \mathbb{1}_{\mathcal{D}_k} ((x_n, t_n)) \end{aligned}$$

8: 以下を満たす T の中で最小の木 T^* を最終的な結果とする.

$$R^{\text{CV}} (T) \leq R^{\text{CV}} (T^{**}) + 1 \times \text{SE} (T^{**})$$

このヒューリスティックな決め方は **1 SE rule** と呼ばれる.

Breiman et al. (1984) の TABLE 3.3

i	\widetilde{T}^i	$R(T^i)$	$R^{CV}(T^i) \pm \text{SE}(T^i)$
1	31	.17	.30 \pm .03
2**	23	.19	.27 \pm .03
3	17	.22	.30 \pm .03
4	15	.23	.30 \pm .03
5	14	.24	.31 \pm .03
6*	10	.29	.30 \pm .03
7	9	.32	.41 \pm .04
8	7	.41	.51 \pm .04
9	6	.46	.53 \pm .04
10	5	.53	.61 \pm .04
11	2	.75	.75 \pm .03
12	1	.86	.86 \pm .03

CARTによる分類

次に分類を扱う。手順は回帰とほぼ同様なので、分割の基準だけ扱う。

$\mathcal{C}_1, \dots, \mathcal{C}_K$ の K クラスに分類する問題を解きたいとする。節点 t を通って流れていくデータの数 $N(t)$ とする。またそれらのうちで \mathcal{C}_k に属するものの数を $N_k(t)$ と表す。このとき

$$p(t|\mathcal{C}_k) = \frac{N_k(t)}{N_k},$$

$$p(\mathcal{C}_k, t) = p(\mathcal{C}_k) p(t|\mathcal{C}_k) = \frac{N_k}{N} \frac{N_k(t)}{N_k} = \frac{N_k(t)}{N},$$

$$p(t) = \sum_{k=1}^K p(\mathcal{C}_k, t) = \frac{N(t)}{N}$$

である。よって事後確率 $p(\mathcal{C}_k|t)$ は次のように求められる。

$$p(\mathcal{C}_k|t) = \frac{p(\mathcal{C}_k, t)}{p(t)} = \frac{N_k(t)}{N(t)}$$

CARTによる分類

分割の基準には**不純度関数** (impurity function) というものを使う. 不純度関数とは, 以下の条件を満たす関数 $\varphi: \Delta^M \rightarrow \mathbb{R}$ のことである.

1. $\varphi(p)$ は $p = (1/M, \dots, 1/M)$ のときに限り最大となる.

2. $\varphi(p)$ は $p = (0, \dots, 0, \overset{i}{\underbrace{1}}, 0, \dots, 0)$ のときに限り最小となる.

3. $\varphi(p_1, \dots, \overset{i}{\underbrace{p_i}}, \dots, \overset{j}{\underbrace{p_j}}, \dots, p_M) = \varphi(p_1, \dots, \overset{i}{\underbrace{p_j}}, \dots, \overset{j}{\underbrace{p_i}}, \dots, p_M).$

不純度

不純度関数を φ とすると、任意の節点 t の不純度 $I(t)$ は

$$I(t) := \varphi(p(\mathcal{C}_1|t), \dots, p(\mathcal{C}_K|t))$$

で定められる.

節点 t の子 t_L と t_R について、 t からそれらに流れるデータの割合を

$$p_L = \frac{p(t_L)}{p(t)}, \quad p_R = \frac{p(t_R)}{p(t)}$$

とすると、 t において基準 s で分割したときの不純度の減少は

$$\Delta I(s, t) = I(t) - p_R I(t_R) - p_L I(t_L)$$

と表される. 不純度の減少が最大になるような分割を行っていけばよい.

不純度関数の例

- ▶ 誤分類率

$$I(t) = 1 - \max_k p(\mathcal{C}_k|t)$$

- ▶ 交差エントロピー

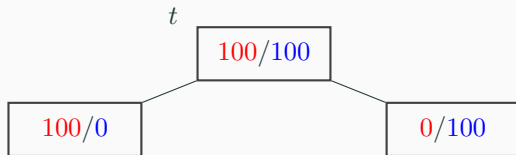
$$I(t) = - \sum_{k=1}^K p(\mathcal{C}_k|t) \log p(\mathcal{C}_k|t)$$

- ▶ ジニ指数 (Gini index)

$$I(t) = \sum_{k=1}^K p(\mathcal{C}_k|t) (1 - p(\mathcal{C}_k|t))$$

不純度の数値例

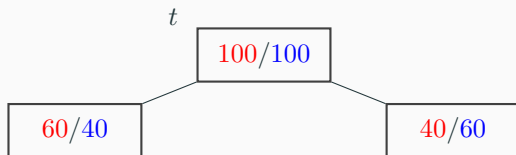
2 クラス分類において、 I として交差エントロピーを使う例.



$$\begin{aligned}\Delta I(s, t) &= I(t) - p_L I(t_L) - p_R I(t_R) \\ &= -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} - 2 \cdot \frac{1}{2} (-1 \log 1 - 0 \log 0) \\ &= \log 2\end{aligned}$$

不純度の数値例

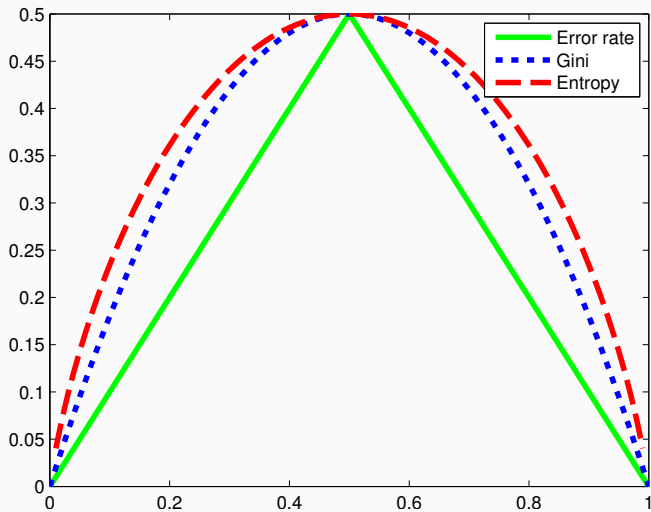
別な分割 s' を使う例.



$$\begin{aligned}\Delta I(s', t) &= I(t) - p_L I(t_L) - p_R I(t_R) \\ &= \log 2 - 2 \cdot \frac{1}{2} \left(-\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} \right) < \Delta I(s, t)\end{aligned}$$

前の分割 s のほうが好ましい.

不純度関数の比較



木構造モデルの問題点

- ▶ 軸に合わせて分割しているので、軸と平行でない方向にうまく境界が引けない.
- ▶ 入力空間をハードに（ただ 1 つのモデルが対応するように）分割してしまう.
- ▶ 通常、回帰では滑らかな関数での近似するが、このモデルでは境界で不連続な、領域ごとの定数予測になってしまう.

参考文献

- Bishop, C. M. (2006). Pattern Recognition and Machine Learning. Springer.
- Breiman, L., Friedman, J., Stone, C. J., and Olshen, R. A. (1984). Classification and regression trees. CRC press.
- Freund, Y. and Schapire, R. E. (1996). Experiments with a new boosting algorithm.
- Friedman, J., Hastie, T., Tibshirani, R., et al. (2000). Additive logistic regression: a statistical view of boosting (with discussion and a rejoinder by the authors). The annals of statistics.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). The Elements of Statistical Learning: data mining, inference and prediction. Springer, second edition.
- Murphy, K. P. (2012). Machine Learning: A Probabilistic Perspective.
- 平井有三 (2012). はじめてのパターン認識. 森北出版.