

Deep Learning 輪読会 2017
第5章 機械学習の基礎

2017.11.13

東京大学大学院経済学研究科 D1 森下光之助
理学系研究科附属天文学教育研究センター B4 吉村勇紀

「Deep Learning」本輪読会 2017

5 章：機械学習の基礎

November 13, 2017

東京大学大学院経済学研究科
澤田研究室 D1 森下光之助

学習アルゴリズム

容量・過剰適合・過少適合

ハイパーパラメータと検証集合

推定量・バイアス・バリエーション

学習アルゴリズム

- 機械学習アルゴリズムとは、データから学習ができるアルゴリズムのこと
- Mitchell (1997) による定義：
「コンピュータプログラムは、性能指標 P で測定されるタスク T における性能が経験 E により改善される場合、そのタスク T のクラスおよび性能指標 P に関して経験 E から学習すると言われている。」

- 機械学習のタスクは機械学習システムがどのように事例 (example) を実行するべきかという観点で記述される
- 事例とは機械学習システムで実行したい対象や事象から定量的に測定された特徴量 (features) の集合である
- 事例はベクトル $\mathbf{x} \in \mathbb{R}^n$ で表す. ここで, ベクトルの各項目 x_i は異なる特徴量である
- 一般的な機械学習のタスクとして, 分類 (Classification), 欠損値のある入力の分類 (Classification with missing inputs), 回帰 (Regression), 転写 (Transcription), 機械翻訳, 構造出力 (Structured output), 異常検知 (Anomaly detection), 合成とサンプリング (Synthesis and sampling), 欠損値補完 (Imputation of missing values), ノイズ除去 (Denoising), 密度推定 (Density estimation) などがある

- 機械学習アルゴリズムの能力を評価するためには、その性能を測る定量的な尺度を設計しなければならない。
 - たとえば分類タスクではモデルの精度（accuracy）を測定することが多い
- 未知のデータに対して機械学習アルゴリズムがうまく機能するかを知りたいので、学習に使われるデータとは異なるテスト集合（test set）を用いて性能指標を評価する

機械学習アルゴリズムは、学習過程においてどのような経験を獲得できるかによって以下の2つに大別できる

- 教師あり学習アルゴリズム (Supervised learning algorithms)
 - 特徴量がラベル (label) や目標 (target) と関連付けられているデータを用いて学習を行う
 - 特徴量 x とラベル y から $p(y|x)$ を予測
 - 分類や回帰など
- 教師なし学習アルゴリズム (Unsupervised learning algorithms)
 - ラベルが付いていないデータ集合から、そのデータ集合構造の有益な特性を学習する
 - 特徴量 x とから $p(x)$ や重要な特性を学習
 - 密度推定やクラスタリングなど

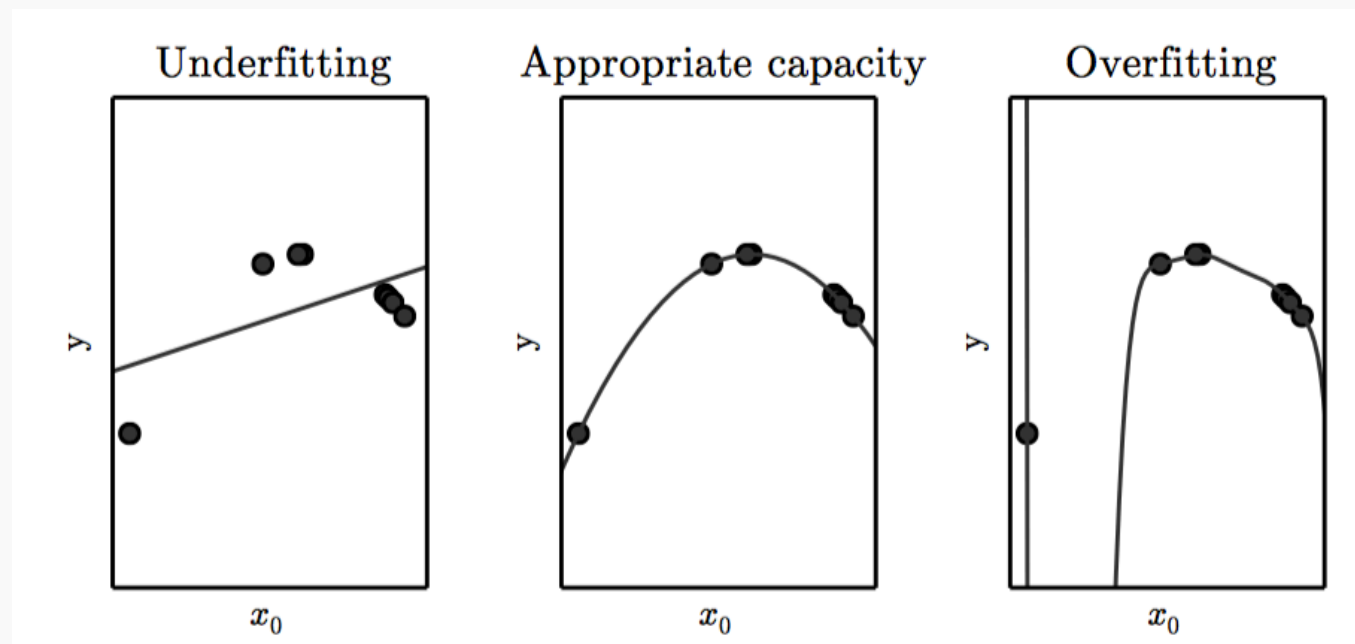
- 大多数の機械学習アルゴリズムでは, 単純にデータ集合を経験する.
(例外: 強化学習 (reinforcement learning))
- 一般的にデータ集合は m 個の要素を含む集合 $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ として記述できる
- 特に $\mathbf{x}^{(i)}$ と $\mathbf{x}^{(j)}$ のサイズが同じ場合はデータ行列 $\mathbf{X} \in \mathbb{R}^{n \times m}$ として記述できる

容量・過剰適合・過少適合

- 機械学習では、モデルの学習に使用した入力だけではなく、これまで見たことのない新たな入力に対しても良い性能を発揮できることが求められる（汎化（generalization））
- 機械学習が最適化と異なるのは、訓練誤差（training error）だけでなく、汎化誤差（generalization error）（テスト誤差（test error）とも呼ばれる）も小さくしたいという点である
- データ生成過程（data-generating process）に、i.i.d. 仮定（i.i.d. assumptions）とを置くことで訓練誤差とテスト誤差の関係性を数学的に調べることが可能になる

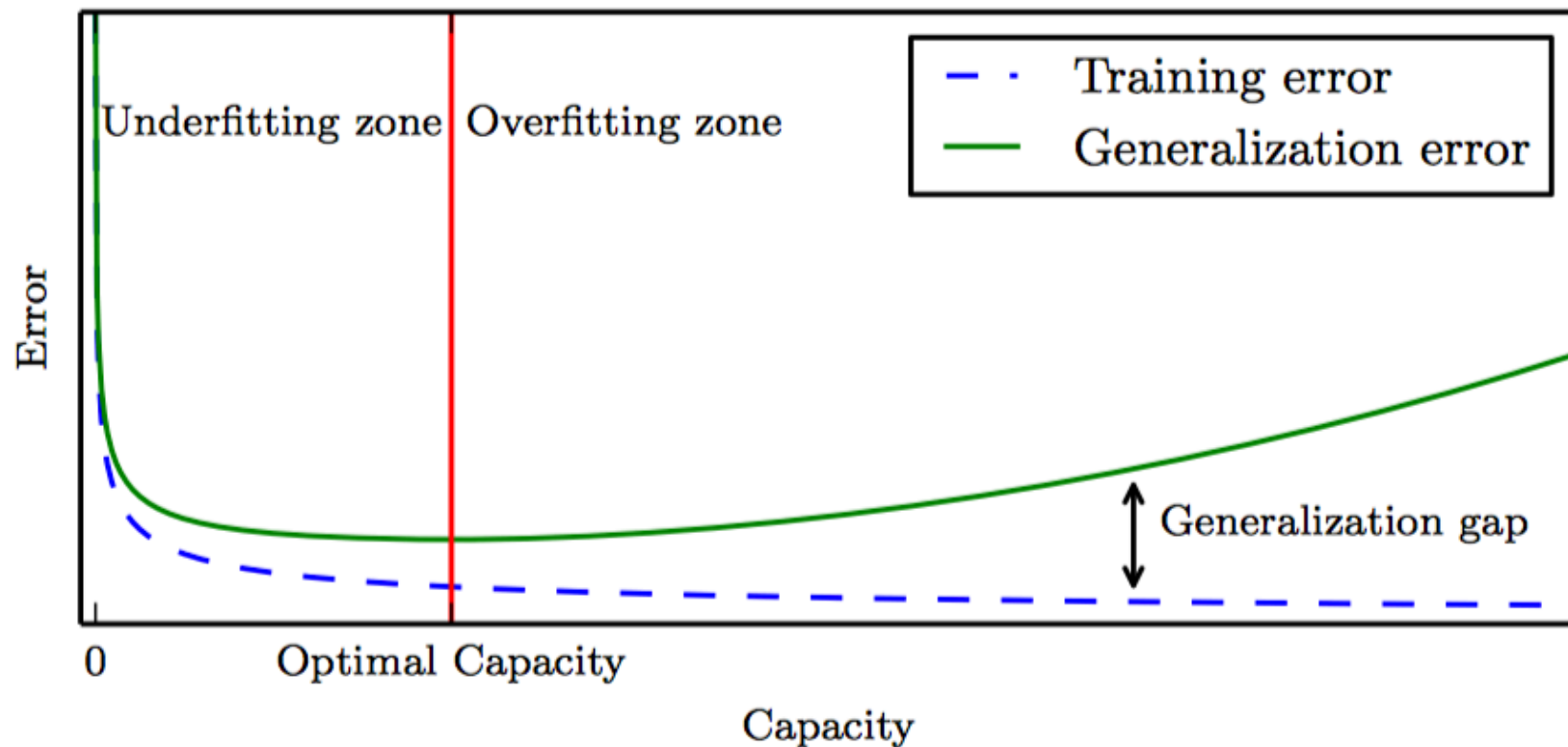
過小適合・過剰適合

- 訓練集合において十分に小さな誤差が得られない場合を過少適合（未学習, underfitting）と呼び、訓練誤差とテスト誤差との差が大きすぎる場合を過剰適合（過学習, overfitting）と呼ぶ
- モデルの容量（capacity）を変更することで過小適合・過剰適合をコントロールする
- 実行する必要のあるタスクの真の複雑さと与えられる訓練データの量に対して適切な容量があるときに、最もよく性能を発揮する



過小適合・過剰適合

- 容量が小さい場合は訓練誤差が大きい（過小適合）
- 容量が増大すると訓練誤差は減少するが、訓練誤差と汎化誤差の差は広がる
- 最終的にはその差が訓練誤差の減少量を上回る（過剰適合）

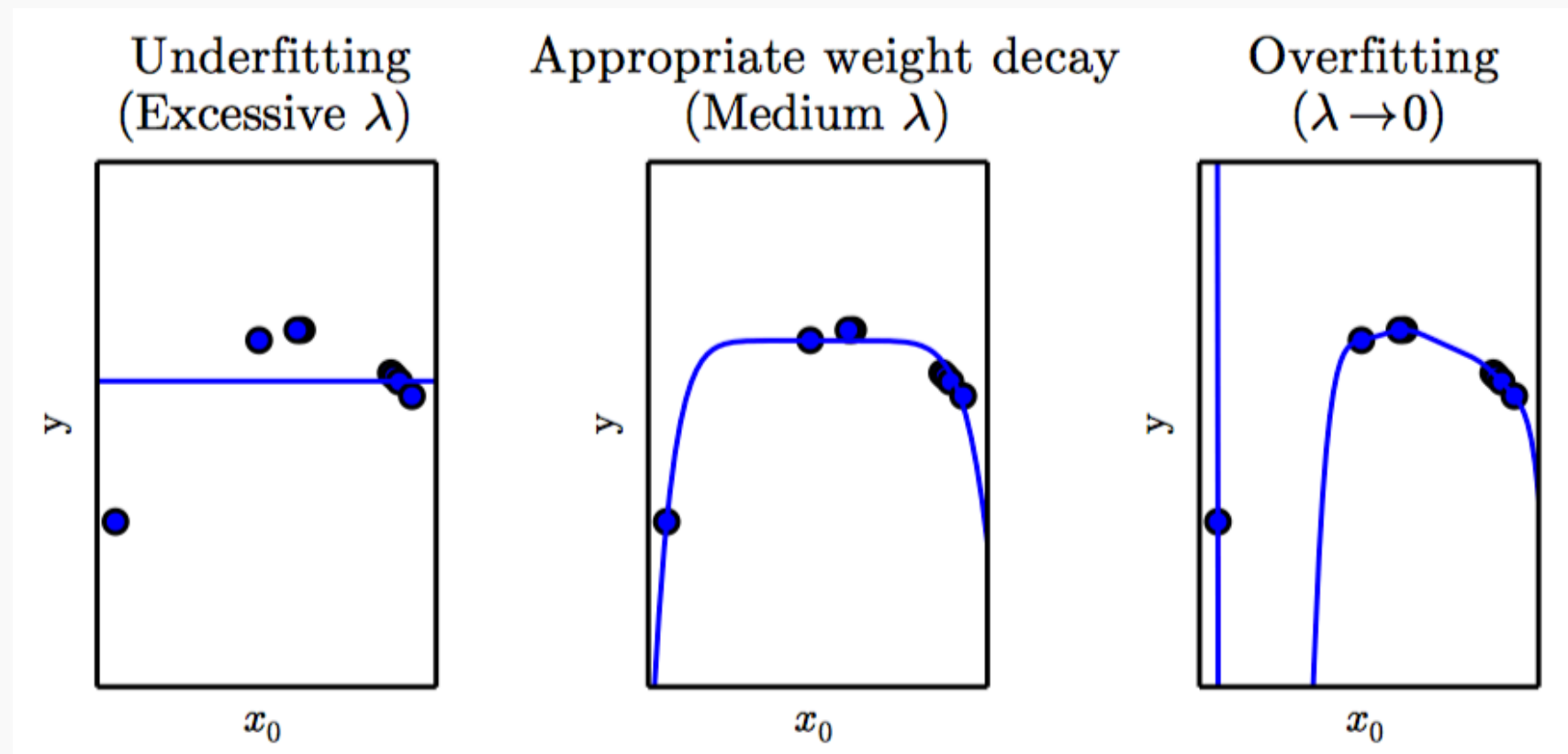


- データを生成する分布すべてを平均すると、どの分類アルゴリズムも、過去に観測されていない点を分類する際の誤差率は同じになる
- つまり、他の機械学習アルゴリズムよりも普遍的に良いと言える機械学習アルゴリズムは存在しない
- よって、普遍的な学習アルゴリズムや絶対的に最良の学習アルゴリズムを探し求めるのではなく、対象としたいデータ 生成の分布から抽出されるデータに対して良い性能を発揮する機械学習アルゴリズムがどのようなものであるかを理解することが目標

正則化

- 正則化項 (regularizer) と呼ばれるペナルティをコスト関数に追加することで, 関数 $f(x; \theta)$ を学習するモデルを正則化できる.
- ex. 重み減衰 (weight decay)

$$J(\mathbf{w}) = \text{MES}_{\text{train}} + \lambda \mathbf{w}^{\top} \mathbf{w}$$



ハイパーパラメータと検証集合

- ほとんどの機械学習アルゴリズムには、そのアルゴリズムの挙動を制御するための設定値がある。この設定値はハイパーパラメータ (hyperparameters) と呼ばれる。
- 訓練集合で学習された場合、このようなハイパーパラメータは常に可能な範囲で最大のモデル容量を選択するので、結果的に過剰適合になる
- この問題を解決するためには、訓練アルゴリズムが観察しない検証集合 (validation set) が必要になる。

k -分割交差検証アルゴリズム

Define $\text{KFoldXV}(\mathbb{D}, A, L, k)$:

Require: 与えられたデータ集合 \mathbb{D} とその要素 $z^{(i)}$

Require: 学習アルゴリズム A , データ集合を入力とし, 学習済みの関数を出力とする関数.

Require: 損失関数 L , 学習済みの関数と事例 $z^{(i)} \in \mathbb{D}$ からスカラー値 $\in \mathbb{R}$ へ変換する関数.

Require: 分割数 k

\mathbb{D} を k 個の互いに排他的な部分テストセット \mathbb{D}_i に分割する. それらの和集合は \mathbb{D} .

for i from 1 to k **do**

$f_i = A(\mathbb{D} \setminus \mathbb{D}_i)$

for $z^{(j)} \in \mathbb{D}_i$ **do**

$e_j = L(f_i, z^{(j)})$

end for

end for

Returne

推定量・バイアス・バリエーション

- 関心のある量（パラメータなど）について「最良の」予測を1つ提示する試みを点推定と呼ぶ
 - パラメータ θ の点推定を $\hat{\theta}$ で表す
- $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ を m 個の独立同一分布 (i.i.d.) のデータ点から成る集合とする. 点推定量 (point estimator) もしくは統計量 (statistic) は, データの任意の関数

$$\hat{\theta}_m = g(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})$$

である

- なお, 現段階では, 真のパラメータの値 θ は固定であるが未知であり, 点推定 $\hat{\theta}$ は確率変数と仮定する

- 良好な推定量は、訓練データを生成した真の潜在的な θ に近い出力を持つ関数
- つまり、バイアスや分散が小さい推定量や、一貫性を持つ推定量はより望ましい推定量と考えられる

- 推定量のバイアスは関数やパラメータの真の値からの期待偏差であり,

$$\text{bias}(\hat{\theta}_m) = \mathbb{E}(\hat{\theta}_m) - \theta_m$$

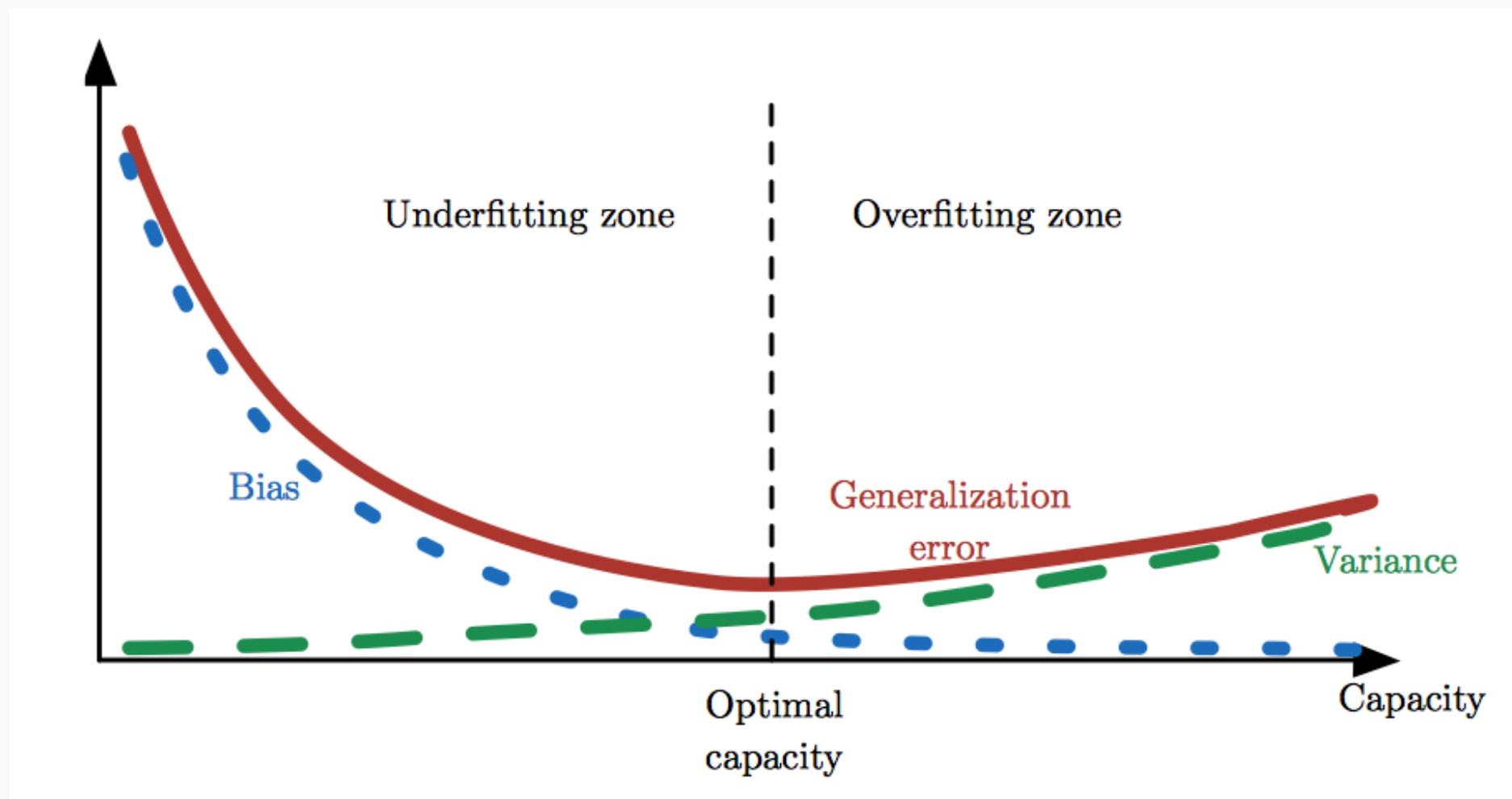
で定義される. バイアスは小さいほうがより良い推定量だと考えられる.

- $\text{bias}(\hat{\theta}_m) = 0$ つまり $\mathbb{E}(\hat{\theta}_m) = \hat{\theta}_m$ のとき推定量は不偏 (unbiased)
- $\lim_{m \rightarrow \infty} \text{bias}(\hat{\theta}_m) = 0$ つまり $\lim_{m \rightarrow \infty} \mathbb{E}(\hat{\theta}_m) = \hat{\theta}_m$ のとき推定量は漸近不偏 (asymptotically unbiased)

- 推定量がデータサンプルの関数としてどれだけ変化すると予測されるかを推定量の分散 (variance) と呼ぶ. 分散の平方根は標準誤差 (standard error) と呼ばれる. 分散が小さいほうがより良い推定量だと考えられる.
- 推定量の分散を $\text{Var}(\hat{\theta})$, 標準誤差を $\text{SE}(\hat{\theta})$ で表す.

バイアスと分散のトレードオフ

一般に、バイアスと分散にはトレードオフの関係がある



推定量の二乗誤差 (mean squared error) を用いることで、バイアスと分散を同時に考慮することが出来る

$$\begin{aligned}\text{MSE}(\hat{\theta}) &= \mathbb{E} \left((\hat{\theta} - \theta)^2 \right) \\ &= \mathbb{E} \left((\hat{\theta} - \mathbb{E}(\theta) + \mathbb{E}(\theta) - \theta)^2 \right) \\ &= \mathbb{E} \left((\hat{\theta} - \mathbb{E}(\theta))^2 \right) + \mathbb{E} \left((\mathbb{E}(\theta) - \theta)^2 \right) \\ &\quad + 2\mathbb{E} \left((\hat{\theta} - \mathbb{E}(\theta)) (\mathbb{E}(\theta) - \theta) \right) \\ &= \text{Var}(\hat{\theta}) + \text{bias}(\hat{\theta})^2\end{aligned}$$

- ここまで固定サイズの訓練集合における推定量の特性に注目してきたが、訓練データの量が増える場合の推定量の挙動にも注意が必要である。
- 特にデータポイント数 m が増加するにつれて点推定量は真の値に収束することが望ましい。これを数式で表すと

$$\text{plim}_{m \rightarrow \infty} \hat{\theta}_m \rightarrow \theta$$

であり、これを（弱）一貫性（consistency）と呼ぶ。ここで plim は確率収束

$$\forall \epsilon > 0, \quad P \left(\left| \hat{\theta}_m - \theta \right| > \epsilon \right) \rightarrow 0 \text{ as } m \rightarrow \infty$$

を意味する。

- なお、 $P \left(\lim_{m \rightarrow \infty} \hat{\theta}_m = \theta \right) = 1$ つまり概収束（Almost sure convergence）する場合には強一貫性をもつという

具体例として、独立同一にガウス分布に従うサンプル集合 $\{x^{(1)}, \dots, x^{(m)}\}$ を考える：

$$\begin{aligned} p\left(x^{(i)}\right) &= \mathcal{N}\left(x^{(i)}; \mu, \sigma^2\right) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\left(x^{(i)} - \mu\right)^2}{2\sigma^2}\right) \end{aligned}$$

- このとき，サンプル平均（sample mean）

$$\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

はガウス平均パラメータ μ の不偏推定量である

$$\mathbb{E}(\hat{\mu}_m) = \mathbb{E}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) = \frac{1}{m} \sum_{i=1}^m \mathbb{E}(x^{(i)}) = \frac{1}{m} \sum_{i=1}^m \mu = \mu$$

- なお，推定量の分散と標準誤差は以下で求まる

$$\text{Var}(\hat{\mu}_m) = \text{Var}\left(\frac{1}{m} \sum_{i=1}^m x^{(i)}\right) = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(x^{(i)}) = \frac{1}{m^2} \sum_{i=1}^m \sigma^2 = \frac{1}{m} \sigma^2$$

$$\text{SE}(\hat{\mu}_m) = \sqrt{\text{Var}(\hat{\mu}_m)} = \frac{\sigma}{\sqrt{m}}$$

- サンプル平均は一致推定量でもある
- 証明にはチェビシェフの不等式

$$\forall \epsilon > 0, \quad P(|X - \mathbb{E}(X)| > \epsilon) < \frac{\text{Var}(X)}{\epsilon^2}$$

を用いる

- $\hat{\mu}_m$ にチェビシェフの不等式を適用すると

$$\forall \epsilon > 0, \quad P(|\hat{\mu}_m - \mu| > \epsilon) < \frac{\sigma^2}{m\epsilon^2} \rightarrow 0 \text{ as } m \rightarrow \infty$$

であり、よって $\hat{\mu}_m$ は一致推定量である

ガウス分布のサンプル分散

- 一方で, サンプル分散 (sample variance)

$$\hat{\sigma}_m^2 = \frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \hat{\mu}_m \right)^2$$

は不偏推定量ではない:

$$\begin{aligned} \mathbb{E}(\hat{\sigma}_m^2) &= \mathbb{E} \left(\frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \hat{\mu}_m \right)^2 \right) \\ &= \mathbb{E} \left(\frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \mu \right)^2 - (\hat{\mu}_m - \mu)^2 \right) \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E} \left(\left(x^{(i)} - \mu \right)^2 \right) - \mathbb{E} \left((\hat{\mu}_m - \mu)^2 \right) \\ &= \frac{m-1}{m} \sigma^2 \end{aligned}$$

- ただし, 漸近不偏ではある:

$$\mathbb{E}(\hat{\sigma}_m^2) = \frac{m-1}{m} \sigma^2 \rightarrow \sigma^2 \text{ as } m \rightarrow \infty$$

- サンプル分散は一致性をもつ

$$\begin{aligned}\hat{\sigma}_m^2 &= \frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \hat{\mu}_m \right)^2 \\ &= \frac{1}{m} \sum_{i=1}^m \left(x^{(i)} - \mu \right)^2 - (\hat{\mu}_m - \mu)^2 \\ &\xrightarrow{p} \sigma^2 \text{ as } m \rightarrow \infty\end{aligned}$$

ガウス分布のサンプル平均（信頼区間）

- $\{x^{(1)}, \dots, x^{(m)}\}$ は独立同一にガウス分布に従うと仮定したので、中心極限定理より

$$\sqrt{m}(\hat{\mu}_m - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2) \text{ as } m \rightarrow \infty$$

- これと連続写像定理を用いると

$$\begin{aligned} \frac{\hat{\mu}_m - \mu}{\text{SE}(\hat{\mu}_m)} &= \frac{\sqrt{m}(\hat{\mu}_m - \mu)}{\hat{\sigma}_m} \\ &\xrightarrow{d} \frac{1}{\sigma} \mathcal{N}(0, \sigma^2) \text{ as } m \rightarrow \infty \\ &= \mathcal{N}(0, 1) \end{aligned}$$

となることがわかる（なお、 \xrightarrow{d} は分布収束を表す）

- よって $\hat{\mu}_m$ の 95%信頼区間は

$$(\hat{\mu}_m - 1.96 \times \text{SE}(\hat{\mu}_m), \hat{\mu}_m + 1.96 \times \text{SE}(\hat{\mu}_m))$$

で与えられる

Deep Learning 輪読会 2017
第5章 機械学習の基礎 (5.5-)

2017.11.13

理学系研究科附属
天文学教育研究センター
学部4年 吉村勇紀

5.5 最尤推定

5.6 ベイズ推定

5.7 教師あり学習アルゴリズム

5.8 教師なし学習アルゴリズム

5.9 確率的勾配降下法

5.10 機械学習アルゴリズムの構築

5.11 深層学習の発展を促す課題

5.5 最尤推定

- 最尤推定
 - モデルを固定したとき、事例集合が発現する確率が最大となるパラメータを求める。

$$\begin{aligned}\boldsymbol{\theta}_{\text{ML}} &= \arg \max_{\boldsymbol{\theta}} p_{\text{model}}(\mathbb{X}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^m p_{\text{model}}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta})\end{aligned}$$

- アンダーフローを考慮して普通 \log をとる

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^m \log p_{\text{model}}(\boldsymbol{x}^{(i)}; \boldsymbol{\theta})$$

5.5 最尤推定

- KLダイバージェンスとの関係
 - 最尤推定は訓練集合で定義される経験分布とモデル分布の差（KLダイバージェンス）を最小化することに相当する

$$D_{\text{KL}}(\hat{p}_{\text{data}} \| p_{\text{model}}) = \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log \hat{p}_{\text{data}}(\mathbf{x}) - \log p_{\text{model}}(\mathbf{x})]$$

- 上式の次の交差エントロピーの最小化と等しい
 - $\mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}} [\log p_{\text{model}}(\mathbf{x})]$
- 最尤推定はモデル分布を経験分布（真のデータ分布は取り扱えない）に一致させる試みである

5.5.1 条件付き対数尤度と平均二乗誤差

- 条件付き対数尤度
 - 入力 x とモデルパラメータに対する出力 y の条件確率から最尤推定量を定式化できる

$$\boldsymbol{\theta}_{\text{ML}} = \arg \max_{\boldsymbol{\theta}} P(\mathbf{Y} \mid \mathbf{X}; \boldsymbol{\theta})$$

- 確率分布がガウシアンで事例が独立同一分布に従う場合

$$\begin{aligned} & \sum_{i=1}^m \log p(y^{(i)} \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= -m \log \sigma - \frac{m}{2} \log(2\pi) - \sum_{i=1}^m \frac{\|\hat{\mathbf{y}}^{(i)} - y^{(i)}\|^2}{2\sigma^2} \end{aligned}$$

5.5.2 最尤法の特徴

- 事例数 $m \rightarrow \infty$ で最尤推定量は真値に漸近する
 - 真の分布がモデル集合内にあり、かつ
 - 真の分布とモデルパラメータが1対1対応している場合
- 最尤推定量の統計的有効性
 - 同じ真値に対する一致推定量でも有限サンプルに対しては汎化誤差が異なる場合がある
 - m が大きい場合、最尤推定量より小さな平均二乗誤差を持つ一致推定量は存在しない

5.6 ベイズ統計

- パラメータの点推定ではなくパラメータの（事後）確率分布を推定する
 - 事後分布 \propto 尤度関数 \times 事前分布

$$p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)}) = \frac{p(x^{(1)}, \dots, x^{(m)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(x^{(1)}, \dots, x^{(m)})}$$

- パラメータの事後分布を畳み込み積分することで最終的な予測を得る

$$p(x^{(m+1)} \mid x^{(1)}, \dots, x^{(m)}) = \int p(x^{(m+1)} \mid \boldsymbol{\theta})p(\boldsymbol{\theta} \mid x^{(1)}, \dots, x^{(m)}) d\boldsymbol{\theta}.$$

5.6 ベイズ統計

- 例) ベイズ線形回帰
 - パラメータの事前分布をガウシアンにとる
 - 途中省略 (本文参照)
 - 事後分布もガウシアンになる
 - 事後分布の平均値を推定量と思うと、二乗和誤差関数に二次正則化項を加えた評価関数で最尤推定した結果と一致

5.6.1 最大事後確率 (MAP) 推定

- MAP推定
 - 事後分布の取り扱いは一般に困難である
 - 点推定量の方が扱いが簡単
 - 事後分布が最大となる値を点推定する
 - ガウシアン事前分布+MAP推定 → 二次正則化
 - 混合ガウス分布などの正則化を設計する際にも持ちいる

5.7 教師あり学習アルゴリズム

- 教師あり学習
 - 入力 x と出力 y の対応を学習する
 - 訓練集合に関してその対応関係は人間の手であらかじめ与えられる（かあるいは自動的に収集される場合もある）

5.7.1 確率的教師あり学習

- 確率的教師あり学習

- 入力に対して線形なモデルを仮定し、出力はモデル値を平均とするガウシアンに従う

$$p(y \mid x; \theta) = \mathcal{N}(y; \theta^\top x, I).$$

- 最適な重みは正規方程式を解くことで一意に定まる。
- クラス分類など2項変数を取り扱う際はロジスティックシグモイド関数などの活性化関数に入れる。

$$p(y = 1 \mid x; \theta) = \sigma(\theta^\top x).$$

- この場合、最適な重みは閉形式の解は存在しない

5.7.2 サポートベクトルマシン

- サポートベクトルマシン
 - 線形関数でモデル化する
 - クラス識別情報のみ出力する
- カーネルトリック
 - モデルは事例との内積で表現できる

- $$w^\top x + b = b + \sum_{i=1}^m \alpha_i x^\top x^{(i)}$$

$$f(x) = b + \sum_i \alpha_i k(x, x^{(i)}).$$

5.7.2 サポートベクトルマシン

- カーネルの利点
 - カーネル関数を計算する方が内積を取るより、計算が簡単
 - カーネルトリックによる改良を施したものを「カーネルマシン」あるいは「カーネル法」と呼ぶ
- カーネルマシンの欠点
 - 決定関数の評価コストが事例数に比例する
 - スパースな α を学習することで評価コストを抑える → サポートベクトル
 - データ集合が大きいと訓練の計算コストが高くなる

5.7.3 その他の教師あり学習アルゴリズム

- k近傍法
 - データ集合のうち入力 x に近い k 個の平均をとる
 - 訓練集合が小さいと汎化性能は悪い
 - 訓練集合が大きくなるほど精度が上がる一方、同時に計算コストは上がる
 - 特徴量間の優劣がつけられない
- 決定木
 - 決定木によって入力空間を（軸に沿って）分割
 - 決定木の訓練は本書の範囲を超える
 - 軸に沿わない決定境界は困難

5.8 教師なし学習アルゴリズム

- 教師なし学習
 - 訓練集合の特徴量を抽出するが教師信号は学習しない
 - データの「最良」の表現を見つける
 - 単純な表現 → 低次元、スパース、独立

5.8.1 主成分分析

- 主成分分析
 - 低次元で相関の小さい表現を学習する
 - 共分散行列が対角行列になるようなデータセットへの線形変換を学習する

$$z = x^{\top} W$$

- 特異値分解（または対角化）によって変換行列を得る

$$X^{\top} X = (U \Sigma W^{\top})^{\top} U \Sigma W^{\top} = W \Sigma^2 W^{\top}$$

5.8.2 k平均クラスタリング

- クラスタリング
 - 入力 x をone-hotコードベクトル h に対応させる
 - i 番目のクラスに属する場合 $h_i=1$ 、それ以外の成分は0
- k平均クラスタリング
 - k 個あるセントロイドの最適化とクラス割り当ての最適化を交互に行い収束させる
- クラスタリングの問題点
 - クラスタリングと現実世界との対応が明らかではない
 - 1つの特徴に対応する複数のクラスタリングが有りうる
 - 分散表現が好まれることもある

5.9 確率的勾配降下法

- 確率的勾配降下法(SGD)

- ほとんどの深層学習はSGDで動作している
- 勾配を用いた最適化

$$\nabla_{\theta} J(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(x^{(i)}, y^{(i)}, \theta)$$

- 計算コストは通常O(m)
- 勾配は期待値なので少ないサンプルだけで評価しても良い
- 最終的に勾配方向にパラメータを進める

$$\theta \leftarrow \theta - \epsilon g$$

5.10 機械学習アルゴリズムの構築

- 機械学習を構成するもの
 - データ集合の仕様
 - 誤差関数
 - 負の対数尤度（+正則化項）
 - モデル
 - 線形 or 非線形
 - 最適化手順

5.11 深層学習の発展を促す課題

- 高次のデータに対する困難
 - 本章で扱う機械学習アルゴリズムでは音声認識や物体認識などの問題を解決できない
 - 高次のデータを扱う際は汎化が指数関数的に困難になる
 - 従来のアルゴリズムは高次データの汎化に適さない

5.11.1 次元の呪い

- 次元の呪い
 - 次元が高いデータには困難が伴う
 - コンポーネントの数 > 事例数
 - 例えば単純な m 次多項式フィッティングだと係数が D^m 個



5.11.2 局所一様と平滑化

- 平滑化事前分布

- 「学習する関数は小さな領域であまり変化してはならない」という仮定を反映した事前分布

$$f^*(\boldsymbol{x}) \approx f^*(\boldsymbol{x} + \epsilon)$$

- k平均法や決定木も多かれ少なかれこの仮定を置いている

- 複雑な関数

- 高次元な関数や領域毎に挙動が違う複雑な関数ではこの仮定は適切でない
- タスク毎に固有の仮定を導入して解決する
- 複数の階層を考えることで仮定を軽くする → ニューラルネット

5.11.3 多様体学習

- 多様体仮説
 - 現実的なデータは高次元空間中の低次元領域に押し込まれている
 - ノイズデータ中に実際のデータを見出す確率は極めて小さい
 - 多様体を描く変換は、実際に想像可能なことが多い
 - 平行移動、回転、ぼかし、変色、…

参考文献

- Deep Learning

- Ian Goodfellow, Yoshua Bengio, Aaron Courville

- 日本語版

<https://www.amazon.co.jp/%E6%B7%B1%E5%B1%A4%E5%AD%A6%E7%BF%92-Ian-Goodfellow/dp/4048930621>