

Wachstumsfunktion

1 Wiederholung: PAC-lernbar (Realisierbarer Fall)

Unsere Aufgabe ist es, Datenpunkte aus einer Menge X zu klassifizieren, beispielsweise $X \subseteq \mathbb{R}$. Die Labels werden binär sein, das heißt -1 oder 1. Beispielsweise könnte X die Menge aller E-Mails sein und die Labels habe die Bedeutung „nicht Spam“ oder „Spam“. Unser Ziel ist es, dass wir für jeden Datenpunkt x , den wir als Eingabe erhalten, das korrekte Label $y \in \{-1, 1\}$ vorhersagen zu können.

Es gibt eine Klasse von Hypothesen \mathcal{H} . Jede hat die Form $h: X \rightarrow \{-1, 1\}$. Wir nehmen an, dass wir im *realisierbaren Fall* sind. Das heißt, es gibt eine Grundwahrheit $f \in \mathcal{H}$, die eine unserer möglichen Hypothesen ist, und das korrekte Label für $x \in X$ ist immer $f(x)$. Wir möchten nun eine Funktion $h \in \mathcal{H}$ finden, die möglichst ähnlich zum korrekten f ist. Dafür steht uns aber nur eine begrenzte Anzahl Samples mit korrekten Labels zur Verfügung.

Wir erinnern uns an die Definition von PAC-Lernbarkeit.

Definition 2.1. Eine Hypothesenklasse \mathcal{H} heißt PAC-lernbar (im realisierbaren Sinn), wenn es eine Funktion $m_{\mathcal{H}}$ und einen Lernalgorithmus \mathcal{A} gibt, sodass der Algorithmus für alle $\epsilon, \delta > 0$, jede Verteilung \mathcal{D} und alle $f \in \mathcal{H}$, gegeben ein Sample S von Größe mindestens $m_{\mathcal{H}}(\epsilon, \delta)$ von Datenpunkten mit korrekten Labels, eine Hypothese $h_S \in \mathcal{H}$ berechnet, sodass $\Pr[\text{err}_{\mathcal{D},f}(h_S) < \epsilon] \geq 1 - \delta$.

Hierbei ist $\text{err}_{\mathcal{D},f}(h) := \Pr_{x \sim \mathcal{D}}[h(x) \neq f(x)]$ der tatsächliche Fehler von h . Zwei Beispiele dafür haben wir bereits gesehen. Heute wollen wir uns das Thema etwas allgemeiner anschauen.

2 Minimierung des Trainingsfehlers

Wir werden uns allgemeiner Algorithmen anschauen, die den Trainingsfehler minimieren.

Definition 2.2. Der Trainingsfehler (oder empirisches Risiko) $\text{err}_S(h)$ einer Hypothese h hinsichtlich einer Trainingsmenge S ist

$$\text{err}_S(h) := \frac{1}{m} |\{h(x_i) \neq y_i\}| .$$

Im realisierbaren Fall gilt für die Grundwahrheit f immer $\text{err}_S(f) = 0$ für alle S . Unsere Algorithmen aus der letzten Vorlesung berechneten jedoch auch jeweils Hypothesen h , sodass $\text{err}_S(h) = 0$. Auch diese minimieren also den Trainingsfehler. Unsere Frage heute wird sein, den tatsächlichen Fehler von Hypothesen zu beschränken, die den Trainingsfehler minimieren.

3 Endliche Hypothesenklassen

Wir betrachten zunächst den einfachen Fall, dass die Menge \mathcal{H} endlich ist, wenn auch ansonsten beliebig.

Satz 2.3. Wenn $m \geq \frac{1}{\epsilon} \ln \left(\frac{|\mathcal{H}|}{\delta} \right)$, dann gilt mit Wahrscheinlichkeit mindestens $1 - \delta$, dass alle $h \in \mathcal{H}$ mit $\text{err}_S(h) = 0$ auch $\text{err}_{\mathcal{D},f}(h) < \epsilon$ erfüllen.

Beweis. Wir betrachten zunächst ein festes $h \in \mathcal{H}$ mit $\text{err}_{\mathcal{D},f}(h) \geq \epsilon$, das heißt, der tatsächliche Fehler von h ist mindestens ϵ . Nun gilt

$$\begin{aligned} \Pr[\text{err}_S(h) = 0] &= \Pr[h(x_1) = y_1, \dots, h(x_m) = y_m] \\ &= \Pr[h(x_1) = y_1] \cdot \dots \cdot \Pr[h(x_m) = y_m] \leq (1 - \epsilon)^m \leq e^{-\epsilon m} . \end{aligned}$$

Das heißt, dass die Wahrscheinlichkeit, dass h keinen Trainingsfehler hat, höchstens $e^{-\epsilon m}$ ist.

Um die Gesamtwahrscheinlichkeit zu beschränken, dass es irgendeine Hypothese gibt, die zwar keinen Trainingsfehler, aber großen tatsächlichen Fehler hat, benutzen wir die sogenannte Union Bound.

Lemma 2.4 (Union Bound). *Es seien $\mathcal{E}_1, \dots, \mathcal{E}_n$ (nicht notwendigerweise disjunkte) Ereignisse. Dann gilt*

$$\Pr\left[\bigcup_{i=1}^n \mathcal{E}_i\right] \leq \sum_{i=1}^n \Pr[\mathcal{E}_i] .$$

Der Beweis der Union Bound folgt durch induktive Anwendung von $\Pr[A \cup B] = \Pr[A] + \Pr[B] - \Pr[A \cap B] \leq \Pr[A] + \Pr[B]$.

Um nun die Union Bound anzuwenden, definieren wir für jede Hypothese $h \in \mathcal{H}$ das Ereignis \mathcal{E}_h , dass $\text{err}_S(h) = 0$.

Nun gilt

$$\begin{aligned} \Pr[\exists h \in \mathcal{H} : \text{err}_{\mathcal{D},f}(h) \geq \epsilon \text{ und } \text{err}_S(h) = 0] &= \Pr\left[\bigcup_{h \in \mathcal{H}, \text{err}_{\mathcal{D},f}(h) \geq \epsilon} \mathcal{E}_h\right] \\ &\leq \sum_{h \in \mathcal{H}, \text{err}_{\mathcal{D},f}(h) \geq \epsilon} \Pr[\text{err}_S(h) = 0] \\ &\leq |\mathcal{H}| e^{-\epsilon m} \leq \delta . \end{aligned} \quad \square$$

4 Wachstumsfunktion

Dieses Ergebnis nützt uns natürlich nichts, wenn \mathcal{H} unendlich ist. Wir haben allerdings schon Beispiele gesehen, dass auch unendliche Hypothesenklassen PAC-lernbar sein können, beispielsweise die Schwellenwertfunktionen. Diese haben eine Struktur, die wir ausnutzen können. Dies können wir wie folgt formalisieren.

Definition 2.5. *Gegeben $S \subseteq X$, sei $\mathcal{H}|_S$ die Menge aller Hypothesen $h \in \mathcal{H}$ mit Definitionsbereich eingeschränkt auf S . Das heißt, $\mathcal{H}|_S = \{h|_S \mid h \in \mathcal{H}\}$.*

Die Wachstumsfunktion von \mathcal{H} ist definiert als $\Pi_{\mathcal{H}}(m) = \max_{S \subseteq X, |S|=m} |\mathcal{H}|_S|$.

Weil die Abbildungen in $\mathcal{H}|_S$ von S nach $\{-1, +1\}$ abbilden, können es nicht mehr als 2^m verschiedene sein, weil es nicht mehr Abbildungen gibt. Somit muss immer $\Pi_{\mathcal{H}}(m) \leq 2^m$ gelten. Häufig sind die Werte von $\Pi_{\mathcal{H}}$ jedoch viel kleiner.

Beispiel 2.6. *Betrachte $X = \mathbb{R}$ und \mathcal{H} als die Klasse der Schwellenwertfunktionen*

$$h_{a'}(x) = \begin{cases} +1 & \text{falls } x \geq a' \\ -1 & \text{sonst} \end{cases}$$

Für $S = \{2, 3, 4\}$ besteht $\mathcal{H}|_S$ aus folgenden vier Funktionen:

$$\begin{array}{ll} x \mapsto -1 & \text{für alle } x \\ x \mapsto \begin{cases} -1 & \text{für } x = 2 \text{ oder } x = 3 \\ +1 & \text{für } x = 4 \end{cases} \end{array} \qquad \begin{array}{ll} x \mapsto +1 & \text{für alle } x \\ x \mapsto \begin{cases} -1 & \text{für } x = 2 \\ +1 & \text{für } x = 3 \text{ oder } x = 4 \end{cases} \end{array}$$

Es gibt noch vier weitere Funktionen $\{2, 3, 4\} \rightarrow \{-1, +1\}$. Diese lassen sich aber nicht über einen Schwellenwert realisieren.

Allgemein gilt $\Pi_{\mathcal{H}}(m) = m + 1$, denn es gibt nur $m + 1$ mögliche „Umschaltunkte“ von -1 auf $+1$. Das heißt, die Funktion wächst deutlich schwächer als 2^m .

Der folgende Satz zeigt, dass wir in der Aussage von Satz 2.3 im Wesentlichen die Größe von \mathcal{H} durch die Wachstumsfunktion ersetzen können.

Satz 2.7. Es seien $\epsilon > 0$ und $\delta > 0$ beliebig und

$$m \geq \max \left\{ \frac{8}{\epsilon}, \frac{2}{\epsilon} \log_2 \left(\frac{2\Pi_{\mathcal{H}}(2m)}{\delta} \right) \right\}. \quad (1)$$

Betrachte ein Sample S von m Datenpunkten mit korrekten Labels gemäß f gezogen unabhängig und identisch verteilt aus \mathcal{D} . Es gilt mit Wahrscheinlichkeit mindestens $1 - \delta$, dass alle $h \in \mathcal{H}$ mit $\text{err}_S(h) = 0$ auch $\text{err}_{\mathcal{D},f}(h) < \epsilon$ erfüllen.

Bevor wir mit dem Beweis dieses Satzes beginnen, schauen wir uns zunächst die Aussage etwas genauer an. Sie hat grundsätzlich die Struktur der Aussage, wie wir sie für PAC-Lernbarkeit brauchen. Wenn m Bedingung (1) erfüllt, dann führt beliebiger Lernalgorithmus, der den Trainingsfehler minimiert, zu einem tatsächlichen Fehler von höchstens ϵ mit Wahrscheinlichkeit mindestens $1 - \delta$.

Wann gilt jedoch Bedingung 1 und wann ist sie überhaupt für alle ϵ und δ erfüllbar? Schauen wir uns nur noch $m \geq \frac{8}{\epsilon}$ an, dann brauchen wir noch

$$m \geq \frac{2}{\epsilon} \log_2 \left(\frac{2\Pi_{\mathcal{H}}(2m)}{\delta} \right) = \frac{2}{\epsilon} \log_2(\Pi_{\mathcal{H}}(2m)) + \frac{2}{\epsilon} \log_2 \left(\frac{2}{\delta} \right) \Leftrightarrow \frac{m - \log_2 \left(\frac{2}{\delta} \right)}{\log_2(\Pi_{\mathcal{H}}(2m))} \geq \frac{2}{\epsilon}.$$

Wenn $\Pi_{\mathcal{H}}(2m) = 2^{2m}$ (die triviale Schranke), dann ist $\log_2(\Pi_{\mathcal{H}}(2m)) = 2m$. Die Ungleichung ist also für sinnvolle ϵ (d.h. $\epsilon < 1$) nicht erfüllbar.

Wächst hingegen $\log_2(\Pi_{\mathcal{H}}(2m))$ schwächer als m , das heißt, $\log_2(\Pi_{\mathcal{H}}(2m)) = o(m)$, dann muss m nur ausreichend groß genug gewählt werden, um die Schranke zu erfüllen.

Im Beispiel mit den Schwellenwertfunktionen ist dies der Fall. Es gilt $\Pi_{\mathcal{H}}(2m) = 2m + 1$. Nun gilt also für alle $\delta > 0$, dass

$$\frac{m - \log_2 \left(\frac{2}{\delta} \right)}{\log_2(\Pi_{\mathcal{H}}(2m))} = \frac{m - \log_2 \left(\frac{2}{\delta} \right)}{\log_2(2m + 1)} \rightarrow \infty \quad \text{für } m \rightarrow \infty.$$

Egal, wie ϵ und δ als gewählt sind, für genügend große m ist Bedingung 1 immer erfüllt.