

Spam-Filterung mit Naive-Bayes-Klassifikator

Team 7: Jannik Weisser & Leon Jerke

Beschreibung des Problems

Im digitalen Zeitalter ist die Flut an E-Mails enorm, und nicht alle davon sind erwünscht oder sogar sicher. Spam-E-Mails, die oft unerwünschte Werbung oder sogar schädliche Inhalte enthalten, sind ein signifikantes Problem für Nutzer und E-Mail-Dienstanbieter. Um diese unerwünschten E-Mails effektiv zu filtern und die Nutzererfahrung zu verbessern, setzen viele E-Mail-Dienste auf maschinelles Lernen, insbesondere auf Naive-Bayes-Klassifikatoren. Diese Klassifikatoren nutzen bedingte Wahrscheinlichkeiten, um zu bestimmen, ob eine eingehende E-Mail Spam ist oder nicht, basierend auf dem Vorkommen bestimmter Schlüsselwörter in der Nachricht.

Formalisierung des Problems

Der Naive-Bayes-Klassifikator basiert auf dem Bayes'schen Theorem, welches bedingte Wahrscheinlichkeiten verwendet, um die Wahrscheinlichkeit einer Hypothese (hier: eine E-Mail ist Spam) auf Basis von vorliegenden Beweisen (hier: Wörter in der E-Mail) zu aktualisieren.

Gegebene Wahrscheinlichkeiten:

1. **Prior-Wahrscheinlichkeit** $P(\text{Spam})$: Die Wahrscheinlichkeit, dass eine beliebige eingehende E-Mail Spam ist, bevor der Inhalt betrachtet wird.
2. **Wahrscheinlichkeiten für das Auftreten bestimmter Wörter in Spam- und Nicht-Spam-E-Mails:**

$P(\text{Wort}|\text{Spam})$: Wahrscheinlichkeit, dass ein bestimmtes Wort in einer Spam-E-Mail erscheint.

$P(\text{Wort}|\neg\text{Spam})$: Wahrscheinlichkeit, dass dasselbe Wort in einer Nicht-Spam-E-Mail erscheint.

Berechnung nach Bayes'schem Theorem:

Das Bayes'sche Theorem wird genutzt, um die Posterior-Wahrscheinlichkeit $P(\text{Spam}|\text{Wörter})$ zu berechnen, die die Wahrscheinlichkeit angibt, dass eine E-Mail Spam ist, gegeben die Wörter, die sie enthält:

$$P(\text{Spam}|\text{Wörter}) = \frac{P(\text{Wörter}|\text{Spam}) \times P(\text{Spam})}{P(\text{Wörter})}$$

Hierbei ist $P(\text{Wörter}|\text{Spam})$ das Produkt der Wahrscheinlichkeiten jedes Wortes, gegeben dass die E-Mail Spam ist, unter der Annahme der Unabhängigkeit der Wörter (daher "naive").

Marginal Likelihood $P(\text{Wörter})$:

Die Gesamtwahrscheinlichkeit $P(\text{Wörter})$, dass die Wörter unabhängig von der Klassifikation auftreten, wird wie folgt berechnet:

$$P(\text{Wörter}) = P(\text{Wörter}|\text{Spam}) \times P(\text{Spam}) + P(\text{Wörter}|\neg\text{Spam}) \times P(\neg\text{Spam})$$

Diese Berechnungen ermöglichen es, jede eingehende E-Mail auf der Basis ihres Inhalts effektiv zu klassifizieren und tragen dazu bei, das Spam-Problem zu minimieren.

Beispiel 1:

Gegebene Daten (fiktiv):

- **Prior-Wahrscheinlichkeit, dass eine zufällige E-Mail Spam ist**
 $P(\text{Spam})$: 20%
- **Prior-Wahrscheinlichkeit, dass eine zufällige E-Mail kein Spam ist**
 $P(\neg\text{Spam})$: 80%

Wir nehmen an, dass bestimmte Schlüsselwörter wie "kostenlos", "Gewinn" und "Sonderangebot" Indikatoren für Spam sein könnten. Wir verwenden folgende fiktive Wahrscheinlichkeiten:

- **Wahrscheinlichkeit des Wortes "kostenlos" in Spam-E-Mails**
 $P(\text{kostenlos}|\text{Spam})$: 30%
- **Wahrscheinlichkeit des Wortes "kostenlos" in Nicht-Spam-E-Mails**
 $P(\text{kostenlos}|\neg\text{Spam})$: 5%

Ziel:

Berechnung der Wahrscheinlichkeit, dass eine E-Mail Spam ist, gegeben das Wort "kostenlos" erscheint.

Formeln:

1. **Bayessches Theorem:**

$$P(\text{Spam}|\text{kostenlos}) = \frac{P(\text{kostenlos}|\text{Spam}) \times P(\text{Spam})}{P(\text{kostenlos})}$$

2. **Marginal Likelihood $P(\text{kostenlos})$:**

$$P(\text{kostenlos}) = P(\text{kostenlos}|\text{Spam}) \times P(\text{Spam}) + P(\text{kostenlos}|\neg\text{Spam}) \times P(\neg\text{Spam})$$

Berechnungen:

1. Zuerst berechnen wir $P(\text{kostenlos})$:

$$P(\text{kostenlos}) = 0.3 \times 0.2 + 0.05 \times 0.8 = 0.1$$

2. Dann nutzen wir das Bayessche Theorem, um $P(\text{Spam}|\text{kostenlos})$ zu berechnen:

$$P(\text{Spam}|\text{kostenlos}) = \frac{0.3 \times 0.2}{0.1} = 0.6$$

Basierend auf den Berechnungen haben wir die folgenden Ergebnisse:

- Die **Gesamtwahrscheinlichkeit, dass das Wort "kostenlos" in einer E-Mail erscheint** ($P(\text{kostenlos})$), beträgt 10%.
- Die **Wahrscheinlichkeit, dass eine E-Mail Spam ist, gegeben dass das Wort "kostenlos" erscheint** ($P(\text{Spam}|\text{kostenlos})$), beträgt 60%.

Beispiel 2:

Gegeben sind die folgenden Emails, die bereits nach Spam oder kein Spam klassifiziert sind:

Nr.	Email	Spam / !Spam
1	send us your password	Spam
2	send us your review	kein Spam
3	review your password	kein Spam
4	review us	Spam
5	send your password	Spam
6	send us your account	Spam

Es soll nun basierend auf diesen Daten, die Wahrscheinlichkeit berechnet werden, dass eine Email mit den Worten "Change your password" Spam ist oder nicht.

Berechnung:

- **Prior-Wahrscheinlichkeit, dass eine zufällige E-Mail Spam ist** $P(\text{Spam}): \frac{4}{6}$
- **Prior-Wahrscheinlichkeit, dass eine zufällige E-Mail kein Spam ist** $P(\neg \text{Spam}): \frac{2}{6}$

Bestimmen der Wortfrequenz:

Nr.	Wort	Wortfrequenz in Spam	Wortfrequenz in \neg Spam
1	password	2/4	1/2
2	review	1/4	2/2
3	send	3/4	1/2
4	us	3/4	1/2
5	your	3/4	2/2
6	account	1/4	0/2

Laplace-Smoothing: Es sollte eine Glättung vorgenommen werden, um das Szenario zu vermeiden, dass ein Wort in den Spam-Trainingsbeispielen nicht vorkommt, wohl aber in den $\neg\text{Spam}$ -Trainingsbeispielen oder umgekehrt. Hierfür wird einfach zu jeder Wortanzahl +1 hinzugefügt und zum Nenner wird +2 hinzugefügt.

Nr.	Wort	Wortfrequenz in Spam	Wortfrequenz in \neg Spam
1	password	3/6	2/4
2	review	2/6	3/4
3	send	4/6	2/4
4	us	4/6	2/4
5	your	4/6	3/4
6	account	2/6	1/4

Als nächstes berechnen wir die Wahrscheinlichkeit dafür, dass "Change your password" in Spam Mails und die Wahrscheinlichkeit dass es in Nicht-Spam Mails vorkommt. Da "Change" nicht in unserem Vokabular vorkommt, kann es einfach weggelassen werden und wir betrachten nur "your password".

- **Spam:**

$$P(\text{your}|\text{Spam}) = \frac{4}{6}$$

$$P(\text{password}|\text{Spam}) = \frac{3}{6}$$

$$P(\text{your password}|\text{Spam}) = P(\text{your}|\text{Spam}) * P(\text{password}|\text{Spam}) = \frac{4}{6} * \frac{3}{6} = \frac{1}{3}$$

- **Nicht Spam:**

$$P(\text{your}|\neg\text{Spam}) = \frac{3}{4}$$

$$P(\text{password}|\neg\text{Spam}) = \frac{2}{4}$$

$$P(\text{your password}|\neg\text{Spam}) = P(\text{your}|\neg\text{Spam}) * P(\text{password}|\neg\text{Spam}) = \frac{3}{4} * \frac{2}{4} = \frac{3}{8}$$

Berechnen der Gesamtwahrscheinlichkeit:

1. Marginal Likelihood:

$$P(\text{your password}) = P(\text{your password}|\text{Spam}) * P(\text{Spam}) + P(\text{your password}|\neg\text{Spam}) * P(\neg\text{Spam}) = \frac{1}{3} * \frac{4}{6} + \frac{3}{8} * \frac{2}{6} = \frac{2}{9} + \frac{1}{8} \approx 0.347$$

2. Bayessches Theorem:

$$P(\text{Spam}|\text{your password}) = \frac{P(\text{your password}|\text{Spam}) * P(\text{Spam})}{P(\text{your password})} = \frac{\frac{1}{3} * \frac{4}{6}}{0.347} = 0.64$$

$$P(\neg\text{Spam}|\text{your password}) = \frac{P(\text{your password}|\neg\text{Spam}) * P(\neg\text{Spam})}{P(\text{your password})} = \frac{\frac{3}{8} * \frac{2}{6}}{0.347} = 0.54$$

Basierend auf den Berechnungen haben wir die folgenden Ergebnisse:

- Die **Gesamtwahrscheinlichkeit, dass die Wörter "your password" in einer E-Mail erscheinen** ($P(\text{your password})$), beträgt 34,7%.
- Die **Wahrscheinlichkeit, dass eine E-Mail Spam ist, gegeben dass die Wörter "your password" erscheinen** ($P(\text{Spam}|\text{your password})$), beträgt 64%.