

به نام خداوند بخشنده مهربان



عنوان تمرین

اسپم یا غیر اسپم

عنوان درس

یادگیری ماشین

استاد

دکتر الهام قصرالدشتی

دستیاران آموزشی

مهرداد قصابی

مریم صفوی

گردآورنده

سید حسین حسینی

بهار ۱۴۰۴

دانشکده مهندسی کامپیوتر

دانشگاه اصفهان

## مستندسازی پروژه: تشخیص ایمیل های هرزنامه (Spam) با طبقه‌بند بیز ساده

### مقدمه

این پروژه با هدف ساخت یک مدل یادگیری ماشین برای تشخیص ایمیل های هرزنامه (Spam) از ایمیل های عادی (Ham) طراحی شده است. تشخیص هرزنامه یکی از کاربردهای کلاسیک و موفق الگوریتم های پردازش زبان طبیعی و طبقه‌بندی متن است. مدل ساخته شده، با تحلیل محتوای یک ایمیل، آن را به یکی از دو دسته "Spam" یا "Ham" طبقه‌بندی می کند.

در این پروژه، همانند پروژه قبلی، از الگوریتم طبقه‌بند بیز ساده (Naive Bayes Classifier) استفاده شده است. این الگوریتم به دلیل سادگی، سرعت بالا و عملکرد بسیار خوب در وظایف طبقه‌بندی متن، انتخاب محبوبی برای این کاربرد است. در اینجا نیز الگوریتم به صورت کامل از پایه و بدون استفاده از کتابخانه های آماده یادگیری ماشین پیاده سازی شده تا مراحل محاسباتی آن به طور شفاف نمایش داده شود. داده های مورد استفاده از فایل spam\_ham\_dataset.csv خوانده شده اند که شامل مجموعه ای از ایمیل ها و برچسب مربوط به آن ها (اسپم یا غیر اسپم) است.

### مراحل اجرای پروژه

پروژه شامل مراحل مختلفی از آماده سازی داده تا ساخت، آموزش و ارزیابی مدل است که در ادامه به تفصیل شرح داده می شوند.

#### ۱. فراخوانی کتابخانه ها (Libraries)

در گام نخست، کتابخانه های ضروری برای انجام پروژه فراخوانی می شوند:

- `seaborn` و `matplotlib.pyplot`: برای مصورسازی داده ها و رسم نمودارها، به ویژه ماتریس درهم ریختگی به صورت نقشه حرارتی. (Heatmap)
- `Pandas`: برای مدیریت و تحلیل داده ها در قالب دیتافریم.
- `Numpy`: برای کار با آرایه های عددی، به خصوص در ساخت ماتریس درهم ریختگی.
- `String`: برای دسترسی به لیستی از علائم نگارشی جهت حذف آن ها از متن.
- `Math`: برای استفاده از تابع لگاریتم در محاسبات احتمالاتی مدل.

• `collections.defaultdict`: برای ساخت دیکشنری‌هایی با مقدار پیش فرض، که در شمارش کلمات بسیار کارآمد است.

## ۲. بارگذاری مجموعه داده (Load Dataset)

مجموعه داده از فایل `spam_ham_dataset.csv` خوانده و در یک دیتافریم `pandas` به نام `df` بارگذاری می‌شود. این دیتافریم شامل ستون‌های اصلی زیر است:

- `text`: محتوای کامل ایمیل که به عنوان ورودی مدل استفاده می‌شود.
- `Label`: برچسب ایمیل که می‌تواند "spam" یا "ham" باشد.
- `label_num`: نسخه عددی برچسب‌ها (معمولاً ۱ برای spam و ۰ برای ham) که برای ارزیابی مدل استفاده می‌شود.

## ۳. پیش‌پردازش متن (Preprocessing Function)

برای آماده‌سازی متن ایمیل‌ها جهت تحلیل، یک تابع به نام `preprocess_text` تعریف شده است. این تابع وظایف زیر را بر عهده دارد:

1. تبدیل به حروف کوچک (Lowercase): یکسان‌سازی حروف برای جلوگیری از تفاوت قائل شدن بین کلماتی مانند "Free" و "free".
2. حذف علائم نگارشی (Punctuation Removal): حذف تمام علائم نگارشی (مانند نقطه، ویرگول، علامت تعجب و غیره) از متن.
3. توکن‌سازی (Tokenization): تقسیم متن پاک‌سازی شده به لیستی از کلمات مجزا (توکن‌ها).

این تابع بر روی ستون `text` دیتافریم اعمال شده و نتایج آن در ستون جدیدی به نام `clean_text` ذخیره می‌شود.

## ۴. تقسیم داده‌ها به دو بخش آموزش و آزمون (Split Train / Test)

برای ارزیابی بی‌طرفانه مدل، داده‌ها به دو بخش تقسیم می‌شوند. در این پروژه، این تقسیم به صورت دستی انجام شده است:

- داده‌های آموزش (Train Data): ۸۰٪ ابتدایی داده‌ها برای آموزش مدل اختصاص داده می‌شود.
- داده‌های آزمون (Test Data): ۲۰٪ انتهایی داده‌ها برای سنجش عملکرد مدل بر روی داده‌های جدید کنار گذاشته می‌شود.

۵. آموزش مدل (ساخت واژگان و شمارش فراوانی‌ها)

در این مرحله، مدل Naive Bayes با استفاده از داده‌های آموزشی، پارامترهای آماری خود را یاد می‌گیرد. این فرآیند شامل مراحل زیر است:

1. ساخت واژگان (Vocabulary): مجموعه‌ای از تمام کلمات منحصر به فرد موجود در ایمیل‌های آموزشی ایجاد می‌شود.

2. شمارش فراوانی کلمات:

○ word\_counts\_spam: تعداد تکرار هر کلمه در ایمیل‌های اسپم.

○ word\_counts\_ham: تعداد تکرار هر کلمه در ایمیل‌های عادی.

3. شمارش اسناد و کلمات:

○ spam\_docs و ham\_docs: تعداد کل ایمیل‌های اسپم و عادی.

○ total\_words\_spam و total\_words\_ham: تعداد کل کلمات در تمام ایمیل‌های اسپم و عادی.

۶. محاسبه احتمالات (Probabilities)

پس از شمارش فراوانی‌ها، احتمالات مورد نیاز برای فرمول Naive Bayes محاسبه می‌شوند:

1. احتمالات پیشین (Prior Probabilities): احتمال اینکه یک ایمیل به صورت تصادفی اسپم یا عادی باشد.

○  $P_{\text{spam}} = \text{spam\_docs} / \text{len}(\text{train\_df})$

○  $P_{\text{ham}} = \text{ham\_docs} / \text{len}(\text{train\_df})$

2. احتمالات شرطی (Conditional Probabilities): احتمال مشاهده یک کلمه خاص، به شرط آنکه

ایمیل از نوع اسپم یا عادی باشد. این محاسبات با استفاده از هموارسازی لاپلاس (Laplace)

(Smoothing) انجام می‌شود تا از بروز احتمال صفر جلوگیری گردد.

۷. پیاده‌سازی طبقه‌بند (Naive Bayes Classifier)

تابع predict بر اساس این احتمالات، یک ایمیل جدید را طبقه‌بندی می‌کند:

1. برای هر ایمیل، امتیاز لگاریتمی تعلق به دسته "spam" و "ham" محاسبه می‌شود.

2. این امتیاز با جمع لگاریتم احتمال پیشین و لگاریتم احتمالات شرطی تمام کلمات موجود در ایمیل به دست می‌آید.

3. در نهایت، دسته‌ای که امتیاز لگاریتمی بالاتری دارد، به عنوان پیش‌بینی نهایی انتخاب می‌شود. (۱ برای اسپم و ۰ برای عادی).

#### ۸. ارزیابی مدل و نمایش نتایج (Evaluation)

عملکرد مدل بر روی مجموعه داده آزمون سنجیده می‌شود:

1. پیش‌بینی بر روی داده‌های آزمون: تابع predict بر روی تمام ایمیل‌های مجموعه آزمون اجرا می‌شود و نتایج در لیست y\_pred ذخیره می‌شود.

2. محاسبه دقت (Accuracy): با مقایسه برچسب‌های واقعی (y\_true) و پیش‌بینی شده (y\_pred)، دقت مدل محاسبه می‌شود. مدل به دقت فوق‌العاده ۹۸.۲۶٪ دست یافته است.

3. ماتریس درهم‌ریختگی (Confusion Matrix): این ماتریس به صورت دستی و سپس با استفاده از Seaborn برای نمایش بهتر، ساخته می‌شود. مقادیر آن عبارتند از:

- مثبت صحیح (TP): تعداد ایمیل‌های اسپم که به درستی اسپم تشخیص داده شده‌اند (۳۰۸ مورد).
- منفی صحیح (TN): تعداد ایمیل‌های عادی که به درستی عادی تشخیص داده شده‌اند (۷۰۹ مورد).

○ مثبت کاذب (FP): تعداد ایمیل‌های عادی که به اشتباه اسپم تشخیص داده شده‌اند (۶ مورد).

○ منفی کاذب (FN): تعداد ایمیل‌های اسپم که به اشتباه عادی تشخیص داده شده‌اند (۱۲ مورد).

این نتایج نشان می‌دهد که مدل در هر دو زمینه تشخیص اسپم و عادی عملکرد بسیار دقیقی دارد و تعداد خطاهای آن بسیار کم است.

#### ۹. ذخیره نتایج (Save Predictions to CSV)

در انتها، پیش‌بینی‌های مدل برای داده‌های آزمون، در یک فایل CSV به

نام spam\_predictions\_from\_scratch.csv ذخیره می‌شود. این فایل شامل یک ستون به

نام Spam است که مقادیر ۰ یا ۱ را در خود دارد.

## جمع‌بندی و نتیجه‌گیری

این پروژه با موفقیت یک سیستم تشخیص هرزنامه را با استفاده از الگوریتم Naive Bayes پیاده‌سازی کرد. دقت بالای ۹۸٪ نشان می‌دهد که این الگوریتم، علی‌رغم سادگی مفروضاتش، برای کاربردهای طبقه‌بندی متن بسیار قدرتمند است. پیاده‌سازی مدل از پایه نیز درک عمیقی از نحوه کارکرد داخلی این الگوریتم فراهم می‌کند.

## منابع و مراجع (References)

1. Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.
  - توضیح: منبعی کلاسیک برای یادگیری مفاهیم Naive Bayes در طبقه‌بندی متن، شامل توضیحات کامل درباره هموارسازی لاپلاس.
2. Jurafsky, D., & Martin, J. H. (2021). *Speech and Language Processing*. 3rd Edition.
  - توضیح: مرجعی جامع در حوزه پردازش زبان طبیعی که به تفصیل به الگوریتم‌های طبقه‌بندی متن می‌پردازد.
3. Scikit-learn Documentation: Naive Bayes.
  - توضیح: برای مقایسه و درک عمیق‌تر انواع مدل‌های Naive Bayes، مستندات کتابخانه Scikit-learn بسیار مفید است.
  - لینک: [https://scikit-learn.org/stable/modules/naive\\_bayes.html](https://scikit-learn.org/stable/modules/naive_bayes.html)
4. Spam/Ham Email Dataset (Kaggle)
  - توضیح: مجموعه داده استفاده شده در این پروژه، یکی از دیتاست‌های استاندارد برای وظایف تشخیص هرزنامه است که در پلتفرم‌هایی مانند Kaggle در دسترس عموم قرار دارد.

😊 موفق باشید 😊