



WIFI LATENCY ANALYSIS REPORT

Aryan Nourbakhsh



1. مقدمه

هدف این پروژه تحلیل داده‌های مربوط به اتصال کاربران به شبکه‌های وای‌فای و بررسی عواملی است که می‌توانند بر مقدار تأخیر (latency) در این شبکه‌ها تأثیرگذار باشند. تأخیر به عنوان یکی از معیارهای اصلی کیفیت سرویس در شبکه‌های بی‌سیم شناخته می‌شود و نقش مهمی در تجربه کاربری به‌ویژه در کاربردهای حساس به زمان مانند تماس تصویری، بازی آنلاین و استریم دارد. در این پروژه سعی شده است با استفاده از تحلیل آماری، مصورسازی داده‌ها و محاسبات اطلاعاتی، متغیرهای تأثیرگذار بر تأخیر شناسایی و میزان اهمیت آن‌ها تعیین شود.

2. معرفی داده‌ها

داده‌های این پروژه در قالب فایل CSV شامل نمونه‌هایی از اتصال کاربران به شبکه‌های وای‌فای در شرایط مختلف می‌باشد. هر ردیف از داده‌ها مربوط به یک اتصال خاص است و ویژگی‌های مختلفی از آن اتصال را ثبت کرده است.

مجموعه داده شامل ویژگی‌های زیر است:

- latency_ms: مقدار تأخیر اتصال به میلی‌ثانیه
- rssi_dbm: قدرت سیگنال دریافتی به دسی‌بل میلی‌وات
- snr_db: نسبت سیگنال به نویز به دسی‌بل
- channel_util%: درصد استفاده از کانال بی‌سیم
- num_assoc_devices: تعداد دستگاه‌های متصل به نقطه دسترسی
- client_speed_mbps: سرعت اتصال کلاینت برحسب مگابیت بر ثانیه
- distance_m: فاصله بین دستگاه کاربر و نقطه دسترسی به متر
- band: باند فرکانسی استفاده‌شده (2.4GHz یا 5GHz)
- Protocol: پروتکل وای‌فای مورد استفاده (مانند 802.11n یا 802.11ac)
- ap_vendor: نام تولیدکننده سخت‌افزار نقطه دسترسی

3. پیش‌پردازش داده‌ها

در این بخش، داده‌ها برای تحلیل آماری و مدل‌سازی آماده‌سازی شدند. مراحل به ترتیب زیر انجام شدند:

3.1. انتخاب ویژگی‌ها

از بین تمامی ستون‌های موجود، تنها ستون‌هایی که به صورت مستقیم بر تأخیر شبکه تأثیرگذار هستند انتخاب شدند تا تمرکز تحلیل روی آن‌ها باشد.

3.2. حذف داده‌های ناقص

تمامی ردیف‌هایی که دارای مقادیر گمشده (NaN) در هر یک از ستون‌های مهم بودند حذف شدند تا از بروز خطا در مراحل بعدی جلوگیری شود.

3.3. حذف نقاط پرت آماری

برای تمامی ویژگی‌های عددی، از روش IQR (Interquartile Range) برای شناسایی و حذف مقادیر پرت استفاده شد. در این روش، بازه‌ای بین $Q1 - 1.5IQR$ و $Q3 + 1.5IQR$ به عنوان بازه مجاز تعریف می‌شود و مقادیری که خارج از این بازه باشند به عنوان پرت حذف می‌شوند.

فکر کنم متوجه نشده باشید. بذارید به توضیح مفصل بهتر بدم.

در تحلیل داده‌ها، نقاط پرت به مقادیری گفته می‌شود که به طور غیرعادی با سایر مقادیر تفاوت دارند. این مقادیر می‌توانند ناشی از خطاهای اندازه‌گیری، شرایط غیرمعمول یا رویدادهای نادر باشند و در صورت عدم حذف، ممکن است باعث تحریف نتایج تحلیل‌های آماری یا مدل‌های پیش‌بینی شوند.

برای شناسایی و حذف این نقاط از روش IQR استفاده می‌شود. IQR یا "فاصله بین چارک‌ها"، تفاوت بین چارک اول ($Q1$) و چارک سوم ($Q3$) داده‌ها است. با محاسبه این فاصله، یک بازه قابل قبول برای داده‌ها تعیین می‌شود که حدود آن برابر است با:

- حد پایین $= Q1 - 1.5 \times IQR$
- حد بالا $= Q3 + 1.5 \times IQR$

مقدارهایی که خارج از این بازه قرار می‌گیرند، به عنوان داده پرت در نظر گرفته شده و از مجموعه داده حذف می‌شوند.

استفاده از روش IQR در مقایسه با روش‌هایی مانند استفاده از میانگین و انحراف معیار، به دلیل عدم حساسیت به توزیع داده و مقاومت در برابر چولگی، روشی مطمئن و مؤثر محسوب می‌شود. به ویژه در داده‌هایی که دارای مقادیر شدید یا غیرنرمال هستند، IQR می‌تواند نقش مهمی در پاک‌سازی داده و بهبود کیفیت تحلیل داشته باشد.

3.4. اعمال فیلترهای منطقی

برخی فیلترها براساس منطق و ماهیت فیزیکی ویژگی‌ها اعمال شد. از جمله:

- فاصله نمی‌تواند منفی باشد
- سرعت اتصال باید بزرگ‌تر یا مساوی صفر باشد
- درصد استفاده از کانال باید در بازه ۰ تا ۱۰۰ باشد
- نسبت سیگنال به نویز باید غیرمنفی باشد

3.5. تبدیل داده‌های دسته‌ای

برای بهینه‌سازی مصرف حافظه و افزایش کارایی پردازش، ستون‌های متنی مانند band ، protocol و ap_vendor به نوع داده category تبدیل شدند.

4. تحلیل ویژگی‌های دسته‌ای

با استفاده از نمودارهای جعبه‌ای (boxplot) ، توزیع تأخیر برای هر مقدار از سه ویژگی دسته‌ای یعنی band ، protocol و ap_vendor مورد بررسی قرار گرفت. این نمودارها نشان دادند که برخی دسته‌ها مانند پروتکل 802.11ac تأخیر کمتری نسبت به دسته‌های قدیمی‌تر دارند. همچنین تفاوت‌هایی بین تولیدکنندگان مختلف نقطه دسترسی مشاهده شد.

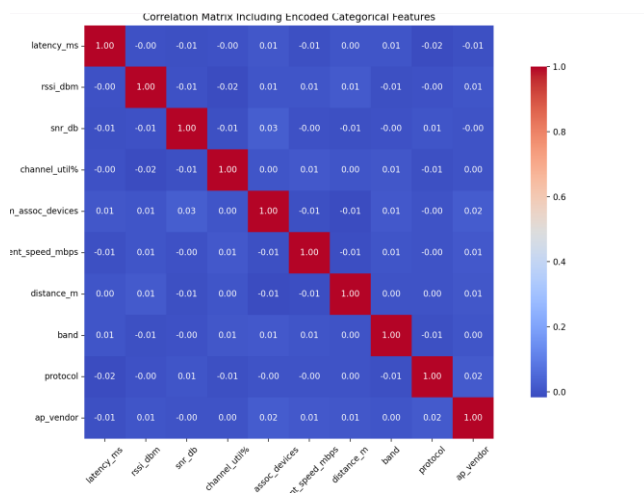
5. تحلیل ویژگی‌های عددی

برای تحلیل روند تأثیر ویژگی‌های عددی مانند فاصله، قدرت سیگنال و استفاده از کانال، هر ویژگی به ۱۰ بازه مساوی تقسیم شد و میانگین تأخیر در هر بازه محاسبه و به صورت نمودار خطی ترسیم شد. نتایج نشان داد که:

- با افزایش فاصله، میانگین تأخیر افزایش می‌یابد.
- با کاهش قدرت سیگنال (RSSI) ، تأخیر بیشتر می‌شود.
- نسبت سیگنال به نویز (SNR) با تأخیر رابطه معکوس دارد.

6. ماتریس همبستگی

برای بررسی همبستگی خطی بین ویژگی‌ها و تأخیر، ابتدا ویژگی‌های دسته‌ای به عدد تبدیل شدند (با استفاده از Label Encoding). سپس ماتریس همبستگی (correlation matrix) محاسبه و به صورت heatmap نمایش داده شد. این تحلیل به طور کلی نشان‌دهنده روابط خطی بین ویژگی‌ها بود ولی نمی‌تواند رابطه‌های غیر خطی را نشان دهد.



7. توزیع لگاریتمی تأخیر

از آن جا که توزیع تأخیر دارای چولگی زیاد است (skewed distribution)، از تبدیل لگاریتمی با تابع \log_{1p} برای نرمال سازی توزیع استفاده شد. سپس هیستوگرام این داده ها ترسیم شد تا توزیع بهتر قابل مشاهده باشد.

8. محاسبه اطلاعات متقابل (Mutual Information)

برای بررسی وابستگی کلی بین ویژگی ها و تأخیر (چه خطی چه غیرخطی)، از روش اطلاعات متقابل استفاده شد. مراحل به شرح زیر است:

8.1. باین بندی داده ها

تمامی ویژگی های عددی با استفاده از KBinsDiscretizer به ۶ بازه مساوی تقسیم شدند. همچنین تأخیر لگاریتمی به ۶ بازه به صورت صدکی (quantile-based) تقسیم شد تا کلاس هدف برای تحلیل اطلاعات متقابل آماده شود. بذارید بیشتر توضیح بدم KBinsDiscretizer یک ابزار از کتابخانه Scikit-learn است که داده های عددی پیوسته را به بازه های گسسته (بندی شده) تبدیل می کند. این ابزار به ویژه زمانی مفید است که بخواهیم ویژگی های پیوسته را برای تحلیل های طبقه بندی یا محاسبه اطلاعات متقابل (Mutual Information) به شکل دسته بندی شده درآوریم.

در این روش، مقدار هر ویژگی عددی به یکی از چند باکس یا بازه مشخص اختصاص داده می شود. تعداد بازه ها توسط کاربر تعیین می شود (مثلاً ۶ بازه)، و می توان نحوه تقسیم بندی را به سه شکل مختلف انجام داد:

- **Uniform**: تقسیم یکنواخت بازه ها بر اساس فاصله عددی
- **Quantile**: تقسیم بر اساس تعداد نمونه ها در هر بازه (هم اندازه از نظر تعداد)
- **Kmeans**: تقسیم با استفاده از خوشه بندی K-Means

در این پروژه از روش uniform استفاده شده که مقدار کل بازه به بازه های مساوی بر اساس فاصله عددی تقسیم می شود.

کاربرد اصلی آن در این پروژه، تبدیل ویژگی های عددی مانند فاصله، قدرت سیگنال و غیره به دسته های گسسته برای امکان مقایسه و محاسبه میزان وابستگی آن ها به تأخیر باین شده (هدف) بوده است. این کار کمک می کند که داده های پیوسته برای تحلیل های آماری یا اطلاعاتی آماده و قابل پردازش شوند.

8.2. رمزگذاری ویژگی های دسته ای

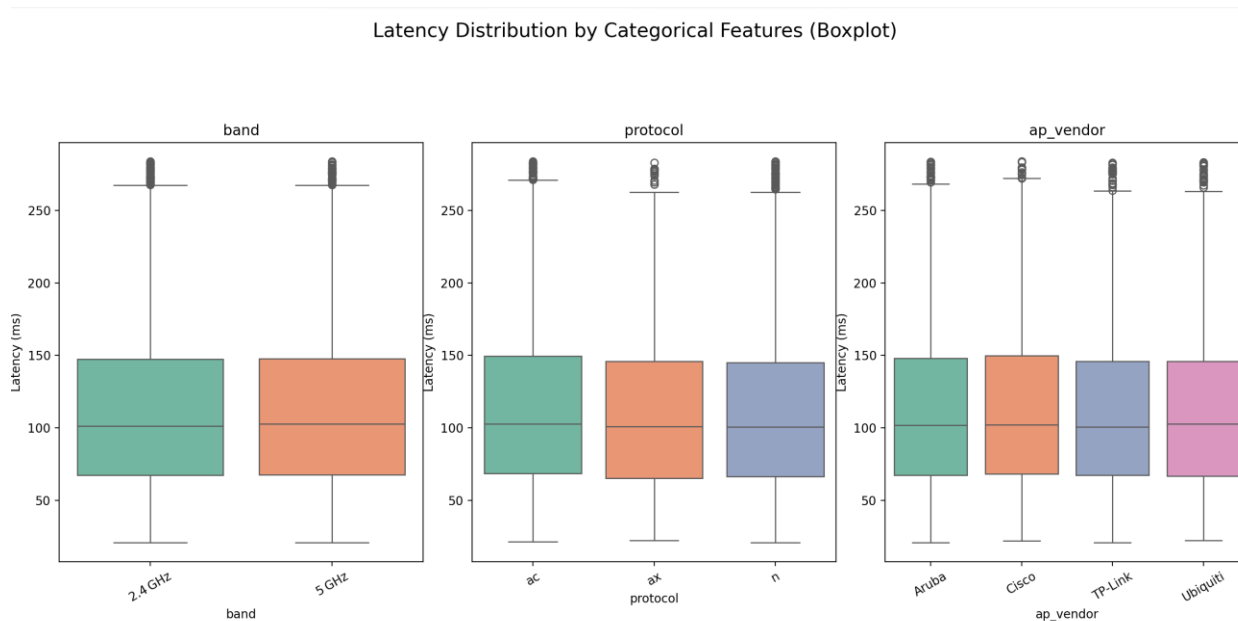
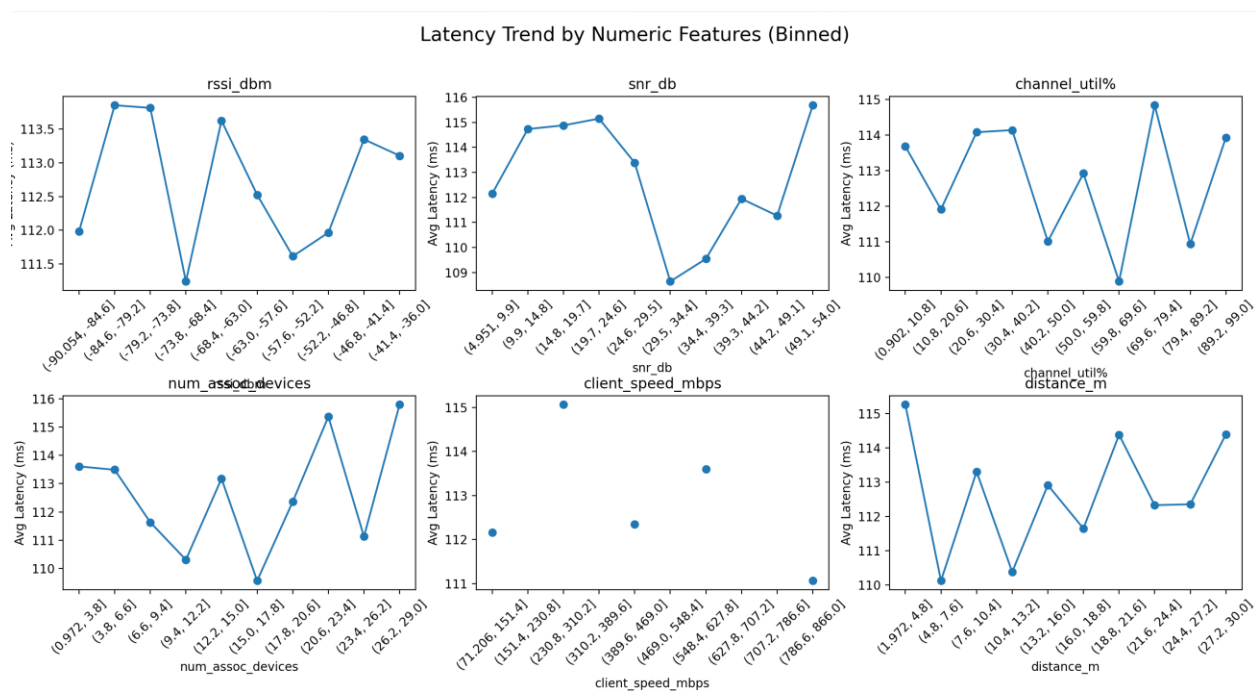
ویژگی های متنی به صورت عددی (Label Encoding) تبدیل شدند.

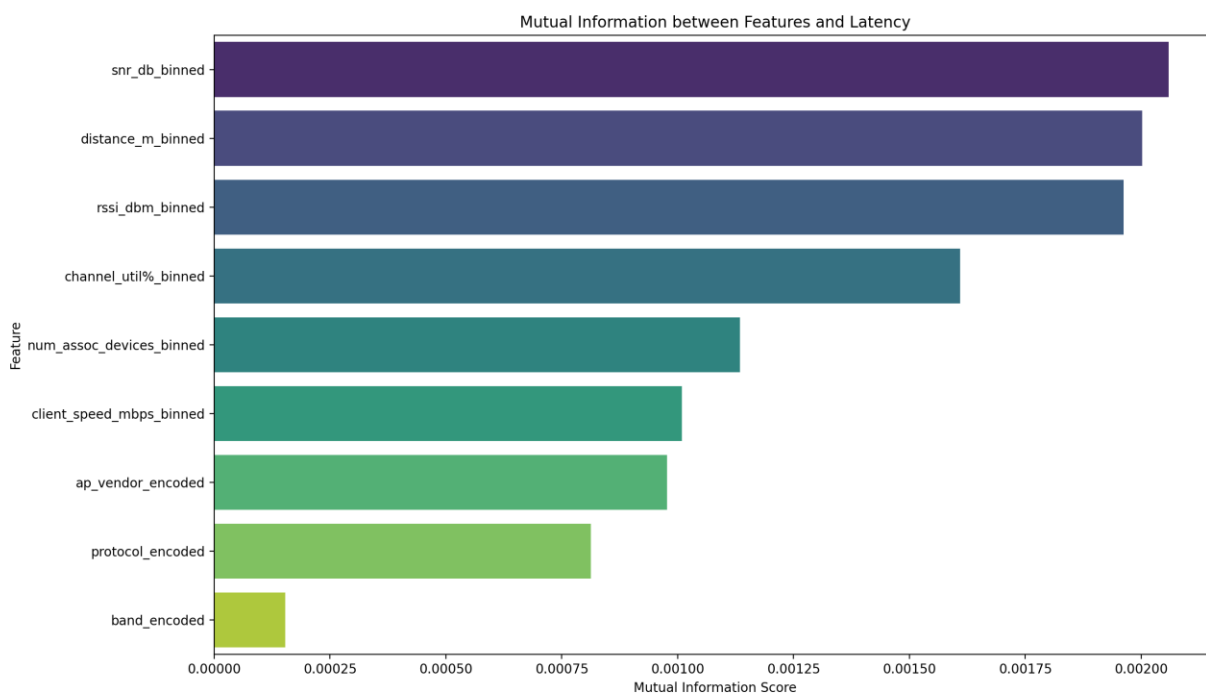
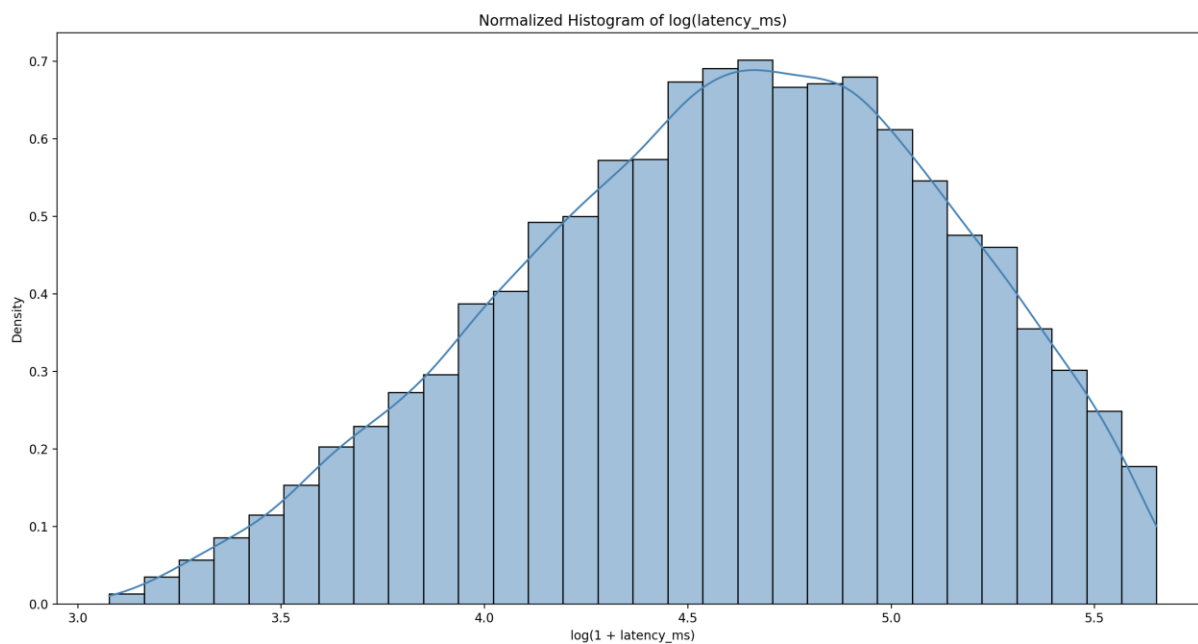
8.3. پیاده سازی تابع اطلاعات متقابل

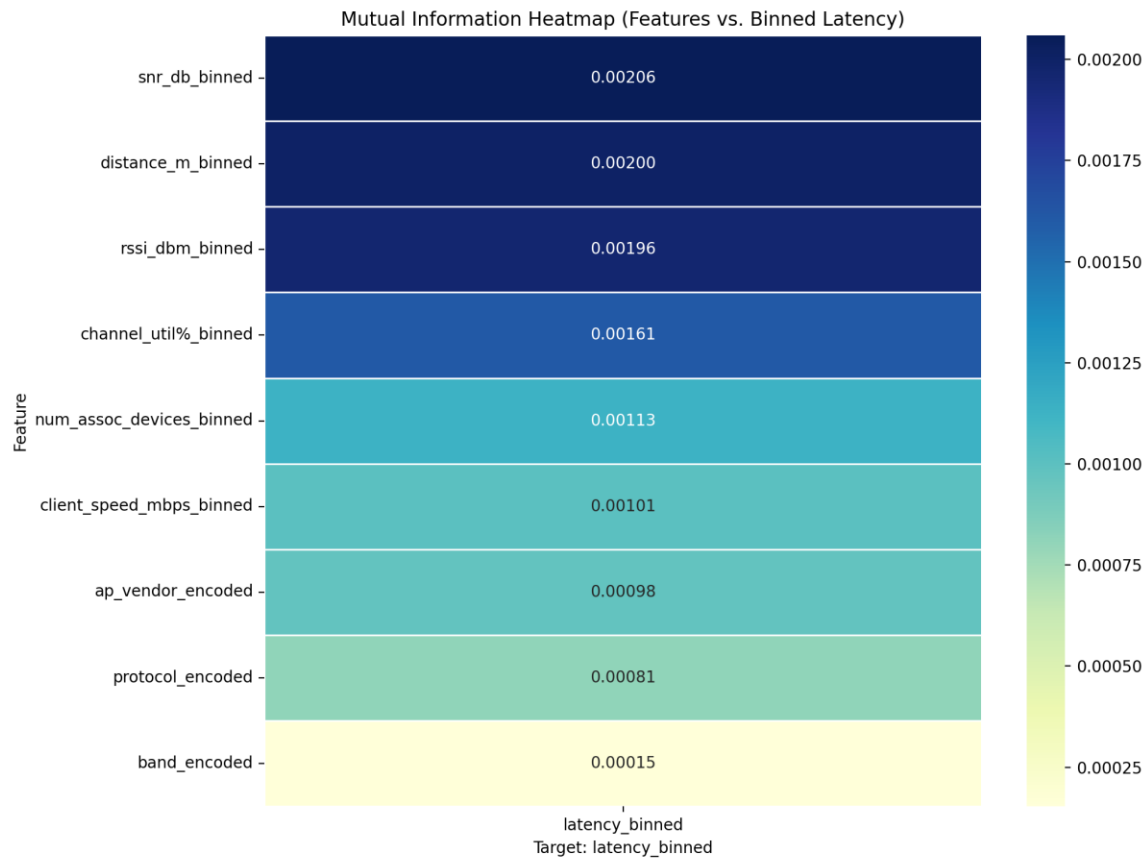
تابعی برای محاسبه اطلاعات متقابل بین دو متغیر به صورت دستی پیاده سازی شد تا بتوان دقیقاً نحوه محاسبه را کنترل کرد و وابستگی بین هر ویژگی و کلاس تأخیر را اندازه گیری نمود.

8.4. ترسیم Heatmap و نمودار میله‌ای

اطلاعات متقابل بین هر ویژگی و تأخیر به صورت heatmap و همچنین نمودار میله‌ای نمایش داده شد. نتایج نشان داد که ویژگی‌هایی مانند فاصله، SNR و RSSI بیشترین اطلاعات را در مورد تأخیر دارند.







بر اساس تحلیل اطلاعات متقابل میان ویژگی‌ها و تأخیر شبکه، نتایج به شرح زیر است:

1. ویژگی‌های `snr_db_binned`، `distance_m_binned` و `rsi_dbm_binned` بیشترین میزان وابستگی به تأخیر دارند و به‌عنوان مهم‌ترین عوامل مؤثر شناسایی شدند. این ویژگی‌ها مستقیماً به کیفیت سیگنال و شرایط فیزیکی ارتباط مربوط می‌شوند.
2. ویژگی‌های `channel_util%_binned`، `num_assoc_devices_binned` و `client_speed_mbps_binned` وابستگی متوسطی به تأخیر دارند. این عوامل بیشتر به وضعیت محیطی و ترافیک شبکه مرتبط هستند.
3. ویژگی‌های `ap_vendor_encoded`، `protocol_encoded` و `band_encoded` کمترین میزان اطلاعات متقابل را با تأخیر دارند. تأثیر این ویژگی‌ها محدودتر و بیشتر غیرمستقیم است.

کد کامل پیاده‌سازی:

<https://github.com/Aryanoor/Wifi-Latency-Analysis.git>