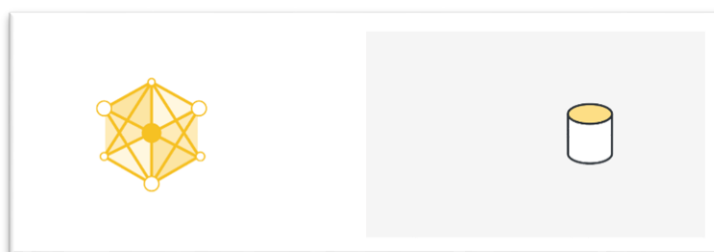


یادگیری فدرال چیست؟

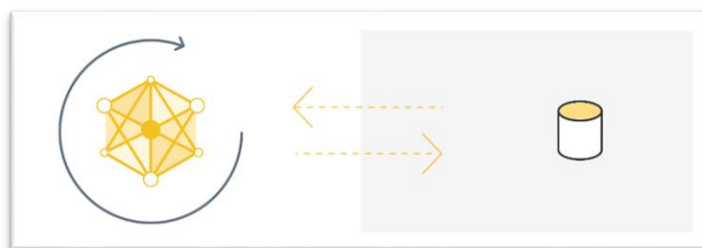
در این گزارش یاد خواهید گرفت که یادگیری فدرال چیست. این آموزش از صفر شروع می شود و انتظار آشنایی با یادگیری فدرال را ندارد. فقط درک پایه ای از علم داده و برنامه نویسی پایتون فرض شده است.

یادگیری ماشین کلاسیک

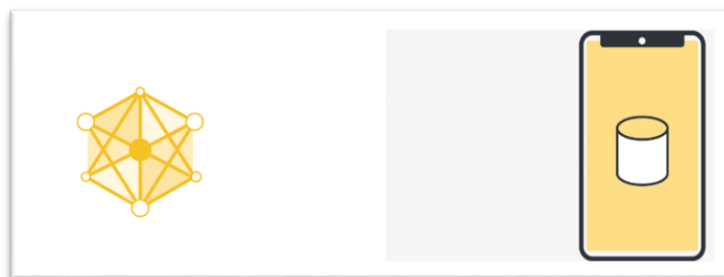
قبل از شروع بحث در مورد یادگیری فدرال، اجازه دهید به سرعت مرور کنیم که امروزه بیشتر یادگیری ماشین چگونه کار می کند. در یادگیری ماشینی، ما یک مدل داریم و داده داریم. مدل می تواند یک شبکه عصبی (همانطور که در اینجا نشان داده شده است) یا چیز دیگری مانند رگرسیون خطی کلاسیک باشد.



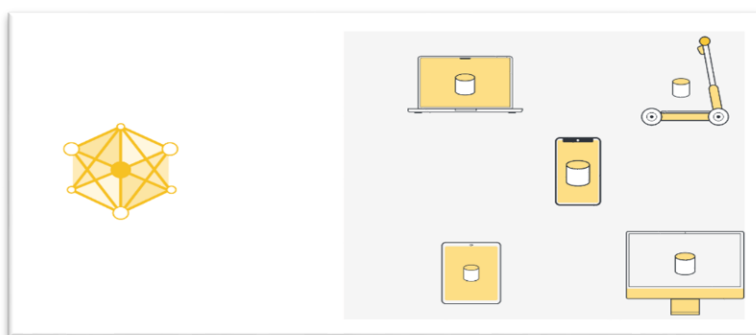
ما مدل را با استفاده از داده ها برای انجام یک کار مفید آموزش می دهیم. یک کار می تواند تشخیص اشیاء در تصاویر، رونویسی یک فایل صوتی ضبط شده یا انجام یک بازی مانند Go باشد.



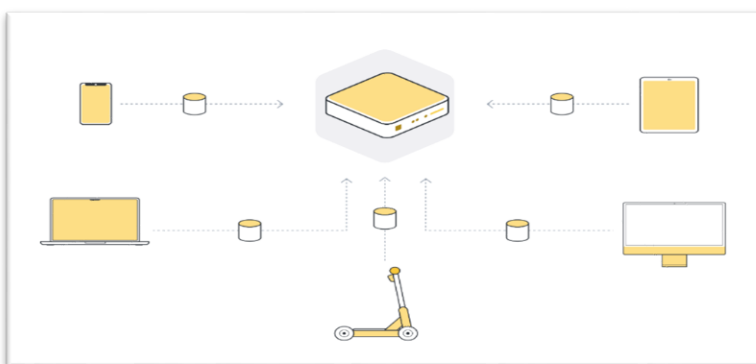
اکنون، در عمل، داده های آموزشی که با آن کار می کنیم، از دستگاهی که مدل را روی آن آموزش می دهیم، منشأ نمی گیرد. در جای دیگری ایجاد می شود. این در گوشی هوشمند توسط کاربر در تعامل با یک برنامه، ماشین جمع آوری داده های حسگر، دریافت ورودی لپ تاپ از طریق صفحه کلید، یا بلندگوی هوشمند در حال گوش دادن به کسی که سعی دارد آهنگی را بخواند، ایجاد می شود.



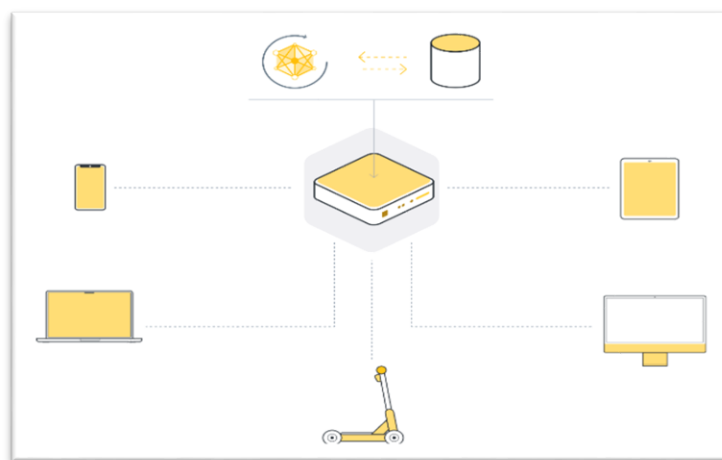
آنچه که ذکر آن نیز مهم است، این "جایی دیگر" معمولاً فقط یک مکان نیست، مکان های زیادی است. ممکن است چندین دستگاه باشد که همه یک برنامه را اجرا می کنند. اما می تواند چندین سازمان نیز باشد که همه داده ها را برای یک کار تولید می کنند.



بنابراین برای استفاده از یادگیری ماشین یا هر نوع تجزیه و تحلیل داده، رویکردی که در گذشته مورد استفاده قرار می گرفت جمع آوری تمام داده ها بر روی یک سرور مرکزی بود. این سرور می تواند جایی در مرکز داده یا جایی در فضای ابری باشد.

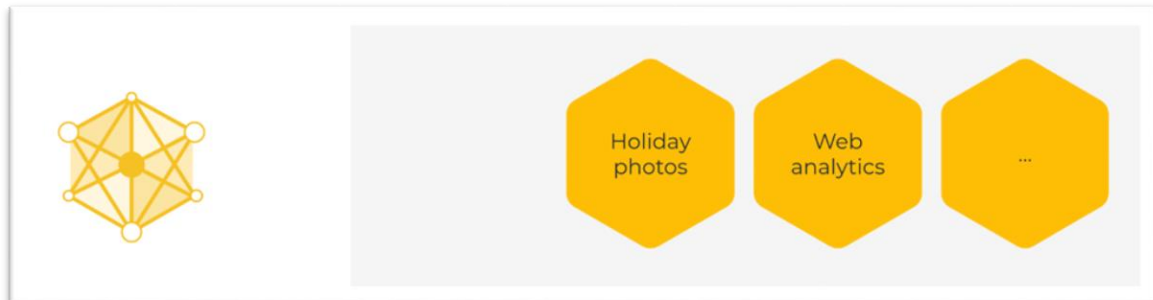


هنگامی که تمام داده ها در یک مکان جمع آوری شد، در نهایت می توانیم از الگوریتم های یادگیری ماشین برای آموزش مدل خود بر روی داده ها استفاده کنیم. این رویکرد یادگیری ماشینی است که ما اساساً همیشه به آن تکیه کرده ایم.

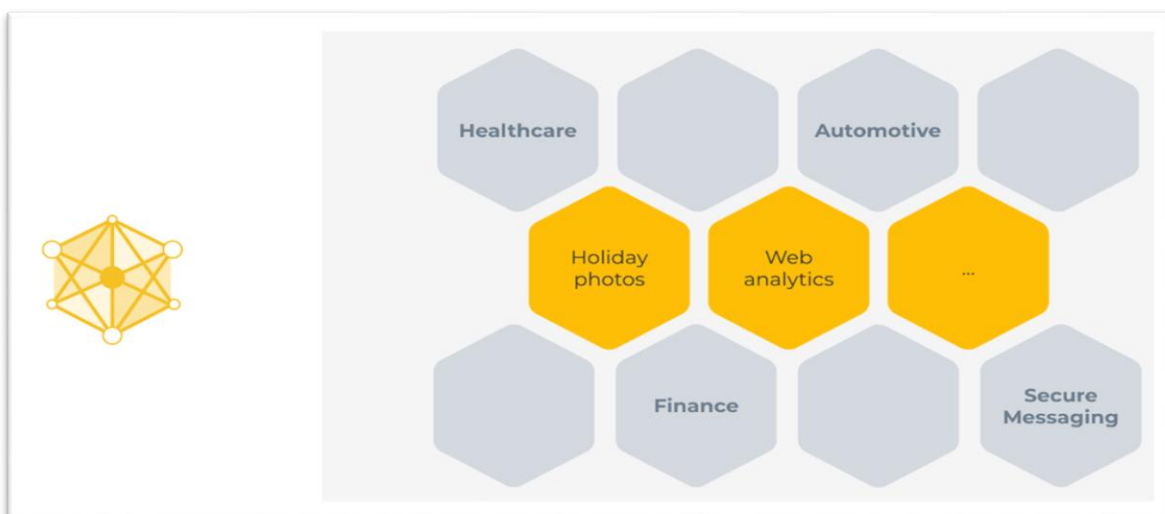


چالش های یادگیری ماشین کلاسیک

رویکرد کلاسیک یادگیری ماشینی که اخیراً دیدیم را می توان در برخی موارد استفاده کرد. نمونه های عالی شامل دسته بندی عکس های تعطیلات، یا تجزیه و تحلیل ترافیک وب است. موارد که در آن تمام داده ها به طور طبیعی در یک سرور متمرکز در دسترس هستند.



اما این رویکرد در بسیاری از موارد دیگر قابل استفاده نیست. مواردی که داده ها در سرور متمرکز در دسترس نیستند یا مواردی که داده های موجود در یک سرور برای آموزش یک مدل خوب کافی نیست.



دلایل زیادی وجود دارد که چرا رویکرد یادگیری ماشین متمرکز کلاسیک برای تعداد زیادی از موارد استفاده بسیار مهم در دنیای واقعی کار نمی کند. آن دلایل عبارتند از:

- **مقررات:** GDPR (اروپا)، CCPA (کالیفرنیا)، PIPEDA (کانادا)، LGPD (برزیل)، PDPL (آرژانتین)، KVKK (ترکیه)، POPI (آفریقای جنوبی)، FSS (روسیه)، CDPR (چین)، PDPB (هند)، PIPA (کره)، APPI (ژاپن)، PDP (اندونزی)، PDPA (سنگاپور)، APP (استرالیا) و سایر مقررات از انتقال داده های حساس محافظت می کنند. در واقع، گاهی اوقات این مقررات حتی از ترکیب کردن داده های کاربران خود برای آموزش هوش مصنوعی توسط سازمان های منفرد جلوگیری می کند، زیرا این کاربران در نقاط مختلف جهان زندگی می کنند و داده های آن ها توسط مقررات حفاظت از داده های متفاوتی کنترل می شود.
- **اولویت کاربر:** علاوه بر مقررات، موارد استفاده ای وجود دارد که کاربران فقط انتظار دارند هیچ داده ای از دستگاه آنها خارج نشود. اگر گذرواژه ها و اطلاعات کارت اعتباری خود را در صفحه کلید دیجیتال گوشی خود تایپ کنید، انتظار ندارید آن رمزهای عبور روی

سرور شرکتی که آن صفحه کلید را توسعه داده است، ختم شود. در واقع، این مورد استفاده دلیل ابداع یادگیری فدرال در وهله اول بود.

- **حجم داده:** برخی از حسگرها، مانند دوربین ها، حجم داده بالایی تولید می کنند که جمع آوری تمام داده ها (به عنوان مثال، به دلیل پهنای باند یا کارایی ارتباط) نه امکان پذیر است و نه اقتصادی. به خدمات ریلی ملی با صدها ایستگاه قطار در سراسر کشور فکر کنید. اگر هر یک از این ایستگاه های قطار مجهز به تعدادی دوربین امنیتی باشد، حجم داده های خام روی دستگاه که تولید می کنند به زیرساخت های فوق العاده قدرتمند و بسیار گرانی برای پردازش و ذخیره سازی نیاز دارد و بیشتر داده ها حتی مفید نیستند.

نمونه هایی که در آن یادگیری ماشین متمرکز کار نمی کند عبارتند از:

- سوابق حساس مراقبت های بهداشتی از چندین بیمارستان برای آموزش مدل های تشخیص سرطان
 - اطلاعات مالی از سازمان های مختلف برای کشف تقلب مالی
 - داده های مکان از ماشین الکتریکی شما برای پیش بینی برد بهتر
 - پیام های رمزگذاری شده سرتاسر برای آموزش مدل های تکمیل خودکار بهتر
- محبوبیت سیستم های افزایش دهنده حریم خصوصی مانند مرورگر Brave یا پیام رسان سیگنال نشان می دهد که کاربران به حفظ حریم خصوصی اهمیت می دهند. در واقع، اگر چنین هشدار وجود داشته باشد، آنها نسخه افزایش دهنده حریم خصوصی را بر سایر گزینه ها انتخاب می کنند. اما برای استفاده از یادگیری ماشین و علم داده در این موارد برای استفاده از داده های خصوصی چه کاری می توانیم انجام دهیم؟ به هر حال، همه اینها حوزه هایی هستند که از پیشرفت های اخیر در هوش مصنوعی به طور قابل توجهی سود می برند.

یادگیری فدرال

یادگیری فدرال به سادگی این رویکرد را معکوس می کند. به جای انتقال داده ها به آموزش، یادگیری ماشین را روی داده های توزیع شده با انتقال آموزش به داده ها امکان پذیر می کند. در اینجا توضیح تک جمله ای آمده است:

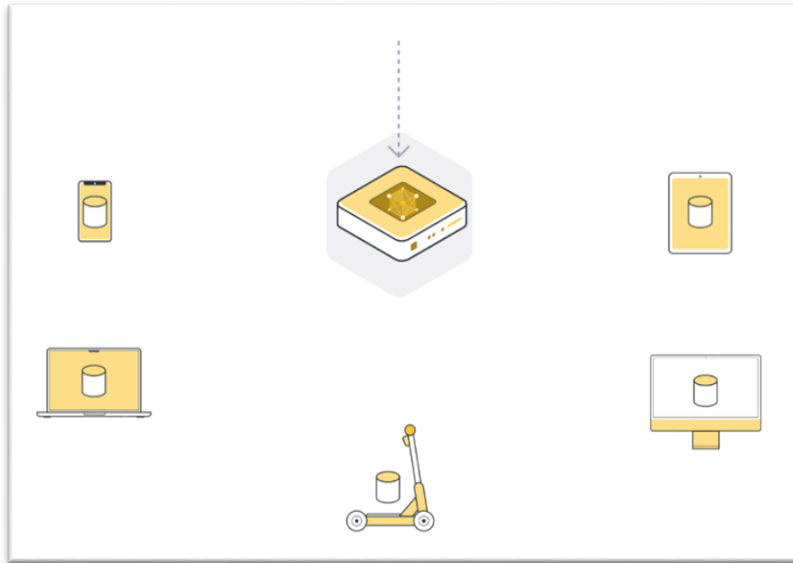
- **یادگیری ماشین مرکزی:** داده ها را به محاسبات منتقل کنید
 - **یادگیری فدرال (ماشین):** محاسبات را به داده ها منتقل کنید
- با انجام این کار، ما را قادر می سازد تا از یادگیری ماشین (و سایر رویکردهای علم داده) در مناطقی استفاده کنیم که قبلاً امکان پذیر نبود. اکنون می توانیم مدل های هوش مصنوعی پزشکی عالی را با ایجاد امکان همکاری در بیمارستان های مختلف آموزش دهیم. ما می توانیم با آموزش مدل های هوش مصنوعی بر روی داده های موسسات مالی مختلف، کلاهبرداری مالی را حل کنیم. ما می توانیم برنامه های جدیدی برای افزایش حریم خصوصی (مانند پیام رسانی امن) بسازیم که هوش مصنوعی داخلی بهتری نسبت به جایگزین های غیرافزاینده حریم خصوصی داشته باشند. و اینها تنها چند نمونه از نمونه هایی است که به ذهن می رسد. همانطور که یادگیری فدرال را به کار می گیریم، مناطق بیشتری را کشف می کنیم که می توانند ناگهان دوباره اختراع شوند، زیرا اکنون به حجم وسیعی از داده های غیرقابل دسترس قبلی دسترسی دارند.

نحوی کار یادگیری فدرال

یادگیری فدرال در ۵ مرحله انجام می گیرد.

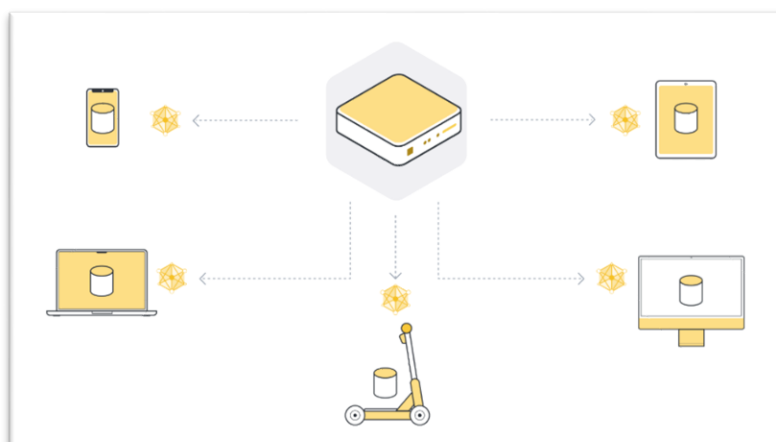
مرحله ۰: مدل جهانی را راه اندازی کنید

ما با مقداردهی اولیه مدل در سرور شروع می کنیم. این دقیقاً در یادگیری متمرکز کلاسیک یکسان است: ما پارامترهای مدل را به صورت تصادفی یا از یک نقطه بازرسی ذخیره شده قبلی مقداردهی اولیه می کنیم.



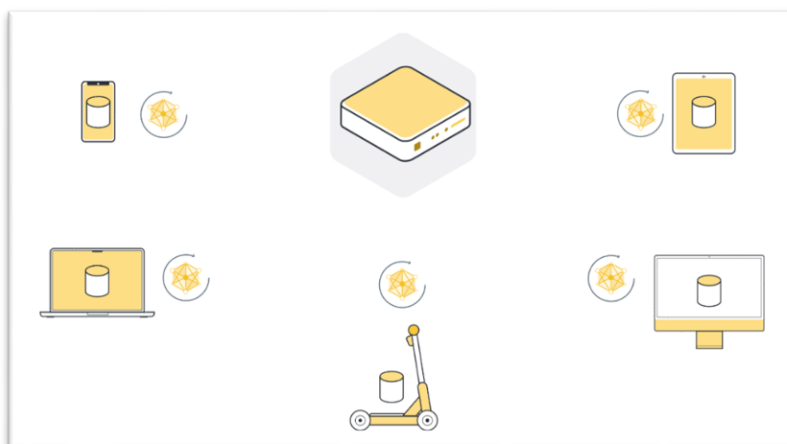
مرحله ۱: ارسال مدل به تعدادی از سازمان ها/دستگاه های متصل (گره های مشتری)

سپس، پارامترهای مدل جهانی را به گره های مشتری متصل می فرستیم (به این فکر کنید: دستگاه های لبه مانند تلفن های هوشمند یا سرورهای متعلق به سازمان ها). این برای اطمینان از این است که هر گره شرکت کننده آموزش محلی خود را با استفاده از پارامترهای مدل مشابه شروع می کند. ما اغلب فقط از تعداد کمی از گره های متصل به جای همه گره ها استفاده می کنیم. دلیل این امر این است که انتخاب بیشتر و بیشتر گره های مشتری بازده کاهشی دارد.



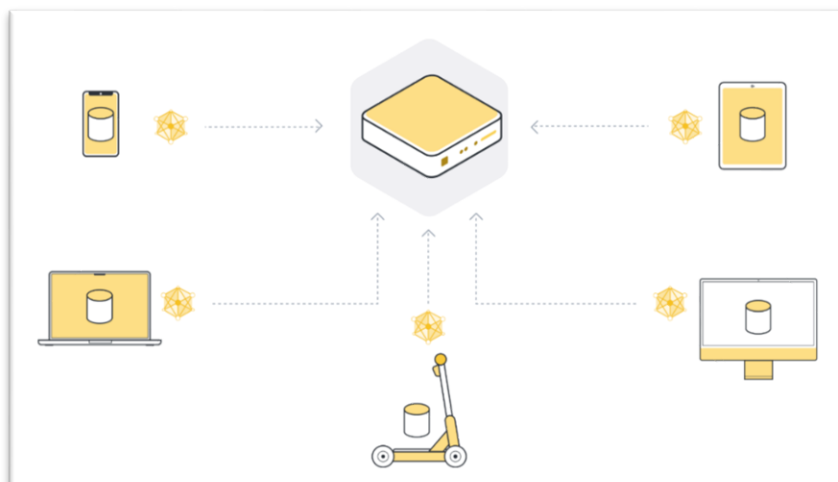
مرحله ۲: اجرا مدل به صورت محلی بر روی داده های هر سازمان/دستگاه (گره مشتری)

اکنون که تمام گره های مشتری (انتخاب شده) آخرین نسخه پارامترهای مدل جهانی را دارند، آموزش محلی را شروع می کنند. آنها از مجموعه داده های محلی خود برای آموزش مدل محلی خود استفاده می کنند. آنها مدل را تا همگرایی کامل آموزش نمی دهند، اما فقط برای مدت کمی تمرین می کنند. این می تواند به اندازه یک دوره در داده های محلی، یا حتی فقط چند مرحله (مینی دسته) باشد.



مرحله ۳: به روز رسانی های مدل را به سرور برگردانید

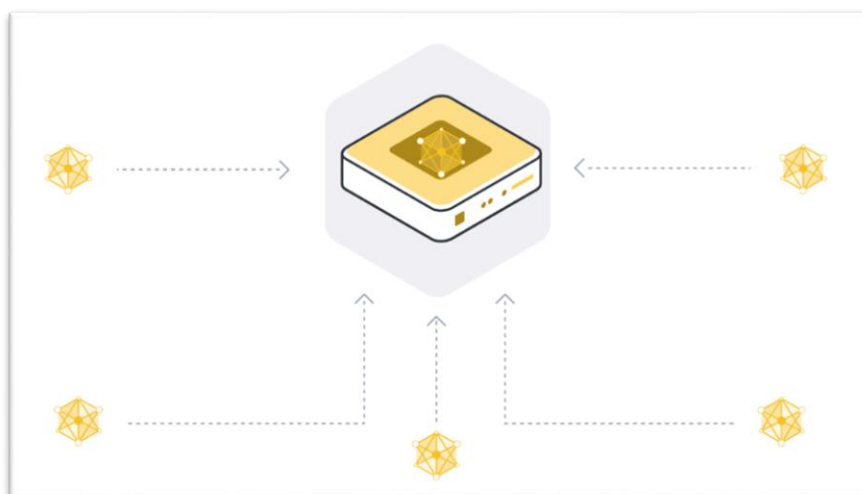
پس از آموزش محلی، هر گره مشتری نسخه کمی متفاوت از پارامترهای مدلی دارد که در ابتدا دریافت کرده بودند. پارامترها همه متفاوت هستند زیرا هر گره مشتری نمونه های متفاوتی در مجموعه داده محلی خود دارد. گره های مشتری سپس آن به روز رسانی های مدل را به سرور ارسال می کنند. به روز رسانی های مدلی که ارسال می کنند می توانند پارامترهای کامل مدل یا فقط گرادیان هایی باشند که در طول آموزش محلی جمع آوری شده اند.



مرحله ۴: به روز رسانی های مدل انبوه در یک مدل جهانی جدید

سرور به روز رسانی های مدل را از گره های مشتری انتخاب شده دریافت می کند. اگر ۱۰۰ گره مشتری را انتخاب کند، اکنون ۱۰۰ نسخه کمی متفاوت از مدل اصلی دارد که هر کدام بر روی داده های محلی یک مشتری آموزش دیده اند. اما آیا ما نمی خواستیم یک مدل داشته باشیم که شامل یادگیری از داده های هر ۱۰۰ گره مشتری باشد؟

برای به دست آوردن یک مدل واحد، باید تمام به روزرسانی های مدلی را که از گره های مشتری دریافت کرده ایم ترکیب کنیم. این فرآیند تجمع نامیده می شود و راه های مختلفی برای انجام آن وجود دارد. اساسی ترین راه برای انجام آن، میانگین گیری فدرال (McMahan et al., ۲۰۱۶) نامیده می شود که اغلب به اختصار FedAvg خوانده می شود. FedAvg ۱۰۰ به روز رسانی مدل را دریافت می کند و همانطور که از نام پیداست، آنها را میانگین می گیرد. برای دقیق تر، میانگین وزنی به روز رسانی های مدل را می گیرد که با تعداد نمونه هایی که هر مشتری برای آموزش استفاده کرده است وزن می شود. وزن دهی برای اطمینان از اینکه هر نمونه داده همان "تأثیر" را بر مدل جهانی حاصل دارد مهم است. اگر یک مشتری ۱۰ نمونه داشته باشد و مشتری دیگر ۱۰۰ مثال داشته باشد، آنگاه - بدون وزن دهی - هر یک از ۱۰ مثال ده برابر هر یک از ۱۰۰ نمونه بر مدل جهانی تأثیر می گذارد.



مرحله ۵: مراحل ۱ تا ۴ را تکرار کنید تا مدل همگرا شود

مراحل ۱ تا ۴ چیزی است که ما آن را دور واحد یادگیری فدرال می نامیم. پارامترهای مدل جهانی به گره های مشتری شرکت کننده ارسال می شوند (مرحله ۱)، گره های مشتری بر روی داده های محلی خود آموزش می بینند (مرحله ۲)، آنها مدل های به روز شده خود را به سرور ارسال می کنند (مرحله ۳)، و سرور سپس به روز رسانی های مدل را جمع می کند. برای دریافت نسخه جدیدی از مدل جهانی (مرحله ۴).

در طول یک دور واحد، هر گره مشتری که در آن تکرار شرکت می کند، فقط برای مدتی تمرین می کند. این بدان معنی است که پس از مرحله تجمع (مرحله ۴)، مدلی داریم که بر روی تمام داده های تمام گره های مشتری شرکت کننده آموزش داده شده است، اما فقط برای مدتی کوتاه. سپس باید این فرآیند آموزشی را بارها و بارها تکرار کنیم تا در نهایت به یک مدل کاملاً آموزش دیده برسیم که در داده های تمام گره های مشتری عملکرد خوبی دارد.

ارزیابی فدرال

همانطور که می توانیم یک مدل را بر روی داده های غیرمتمرکز گره های مشتری مختلف آموزش دهیم، می توانیم مدل را روی آن داده ها برای دریافت معیارهای ارزشمند ارزیابی کنیم. این ارزیابی فدرال نامیده می شود که گاهی اوقات به اختصار FE خوانده می شود. در واقع، ارزیابی فدرال بخشی جدایی ناپذیر از اکثر سیستم های یادگیری فدرال است.

تجزیه و تحلیل فدرال

در بسیاری از موارد، یادگیری ماشینی برای استخراج ارزش از داده‌ها ضروری نیست. تجزیه و تحلیل داده‌ها می‌تواند بینش‌های ارزشمندی به دست آورد، اما باز هم، اغلب داده‌های کافی برای دریافت پاسخ روشن وجود ندارد. میانگین سنی که در آن افراد دچار یک نوع بیماری خاص می‌شوند چقدر است؟ تجزیه و تحلیل فدرال چنین پرس و جوایی را در چندین گره مشتری فعال می‌کند. معمولاً در ارتباط با سایر فناوری‌های تقویت‌کننده حریم خصوصی مانند تجمیع امن استفاده می‌شود تا از مشاهده نتایج ارسال شده توسط گره‌های مشتری منفرد توسط سرور جلوگیری شود.

حریم خصوصی دیفرانسیل

حریم خصوصی متفاوت (DP) اغلب در زمینه یادگیری فدرال ذکر می‌شود. این یک روش حفظ حریم خصوصی است که هنگام تجزیه و تحلیل و به اشتراک گذاری داده‌های آماری استفاده می‌شود و از حریم خصوصی افراد شرکت‌کننده اطمینان می‌یابد. DP با افزودن نویز آماری به به‌روزرسانی‌های مدل، به این مهم دست می‌یابد و اطمینان می‌دهد که اطلاعات هر شرکت‌کننده‌ای قابل تشخیص یا شناسایی مجدد نیست. این تکنیک را می‌توان یک بهینه‌سازی در نظر گرفت که یک معیار حفاظت از حریم خصوصی قابل سنجش را ارائه می‌دهد.

کل

یادگیری فدرال، ارزیابی فدرال و تجزیه و تحلیل فدرال به زیرساخت نیاز دارند تا مدل‌های یادگیری ماشینی را به جلو و عقب منتقل کنند، آن‌ها را بر اساس داده‌های محلی آموزش و ارزیابی کنند، و سپس مدل‌های به‌روز شده را تجمیع کنند. Flower زیرساختی را برای انجام دقیقاً به روشی آسان، مقیاس پذیر و ایمن فراهم می‌کند. به طور خلاصه، Flower یک رویکرد یکپارچه برای یادگیری، تجزیه و تحلیل و ارزشیابی فدرال ارائه می‌دهد. این به کاربر اجازه می‌دهد تا هر حجم کاری، هر چارچوب ML و هر زبان برنامه نویسی را فدرال کند.

