

Gaussian discriminant analysis

مثالی که جلسه قبل بررسی کردیم یک مثال از یک دسته از مدل ها بود به نام GDA . تو GDA ما فرضی که میکنیم که توزیع های $conditional$ ما به سری توزیع نرمال هستند و اینکه $prior$ روی کلاس های ما برنولی هست. (فرض ما حالت دو کلاس مساله است).

$$\mathcal{C}_k \sim \text{Bernouli}(\phi) = \phi^y(1 - \phi)^{(1-y)}$$

$$p(\mathbf{x}|y = k) \sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

$$P(X = x|Y = 1) = \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}_1|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \boldsymbol{\Sigma}_1^{-1}(x - \mu_1)\right)$$

$$P(X = x|Y = 0) = \frac{1}{(2\pi)^{\frac{n}{2}} |\boldsymbol{\Sigma}_0|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \boldsymbol{\Sigma}_0^{-1}(x - \mu_0)\right)$$

ما تا الان $prior$ کلاس ها را به صورت فرض شده در نظر گرفتیم ولی الان کمی مدل بهتری میتونیم اتخاذ کنیم. یک تابع $likelihood$ تشکیل میدیم با فرض iid بودن سمپل ها:

$$\begin{aligned} \ell(\phi, \mu_0, \mu_1, \Sigma_1, \Sigma_2) &= \log \prod_{i=1}^N p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma_0, \Sigma_1) \\ &= \log \prod_{i=1}^N p(x^{(i)}|y^{(i)}, \mu_0, \mu_1, \Sigma_0, \Sigma_1) p(y^{(i)}; \phi) \end{aligned}$$

با ماکزیمم کردن ℓ با توجه به پارامترها، تخمین $likelihood$ پارامترها را پیدا می کنیم.

- برای به حداکثر رساندن $\ell(\cdot)$ با توجه به پارامترها، مشتق $\ell(\cdot)$ را می‌گیریم، مشتق را برابر ۰ قرار می‌دهیم و سپس مقادیر پارامترهایی که عبارت $\ell(\cdot)$ را به حداکثر می‌رساند را حل می‌کنیم. این تخمین حداکثر احتمال ϕ را به دست می‌دهد.

توضیح شهودی در مورد مقدار ϕ برای به حداکثر رساندن *Likelihood* به شرح زیر است. به یاد بیاورید که ϕ تخمین احتمال y برابر با ۱ است. در مثال خاص ما، این شانس که بیمار بعدی با تومور بدخیم به مطب پزشک شما می‌رود با ϕ نشان داده می‌شود. احتمال شیر آمدن در پرتاب سکه، کسری از پرتاب‌های دیده شده است. به همین ترتیب، برآورد حداکثر احتمال برای ϕ فقط کسری از نمونه‌های آموزشی شما با برچسب $y = 1$ است.

$$\begin{aligned}\phi &= \frac{1}{n} \sum_{i=1}^n 1\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}\end{aligned}$$

- برای ایجاد شهود پیرامون مقدار μ_0 و μ_1 که مقدار فرمول را به حداکثر می‌رساند، به این فکر کنید که تخمین حداکثر *likelihood* میانگین همه ویژگی‌ها برای یک کلاس خاص، مثلاً تومورهای خوش خیم، که با برچسب کلاس ۰ در مجموعه داده ما نشان داده می‌شود، چقدر خواهد بود.

μ_0 را در نظر بگیرید. یک راه معقول برای تخمین μ_0 این است که تمام تومورهای خوش خیم *training set* خود را در نظر بگیرید (همه نمونه‌های منفی، یعنی ورودی‌ها به صورت ۰) و فقط میانگین ویژگی‌های آنها را در نظر بگیرید. معادله فوق راهی برای نوشتن این شهود است. صورت کسر مجموع بردارهای ویژگی برای همه نمونه‌های تومورهای خوش خیم در *training set* است (یعنی نمونه‌هایی با $y = 0$) در حالی که مخرج به سادگی تعداد تومورهای خوش خیم در مجموعه آموزشی است.

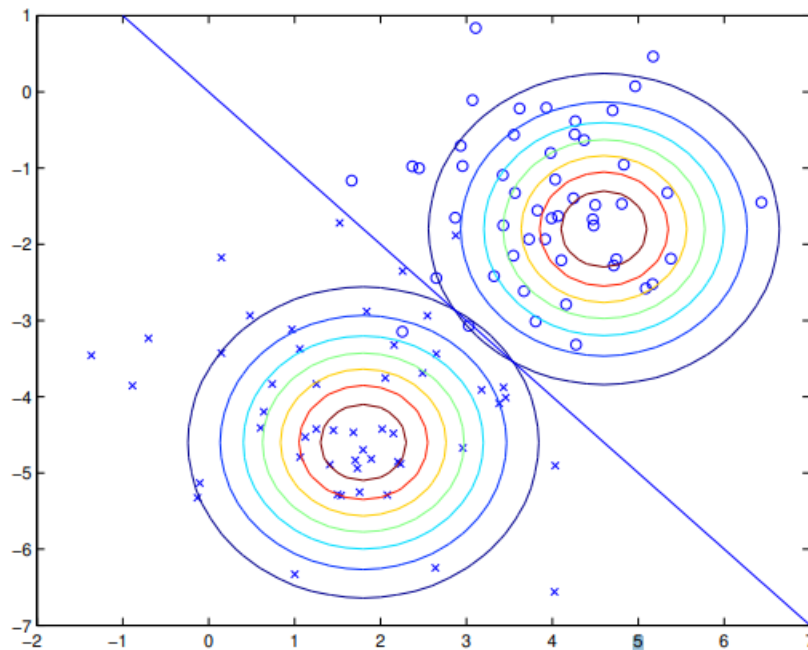
البته بدیهه‌ها وقتی y را میدانیم (از *observation*) با توجه به ستینگ احتمالاتی ما دیگه پارامترهای دیگر توضیح رو نیاز نداریم و توضیح پارامتر *class conditional* هم وجودش معنایی ندارد.

مورد خاص: $\Sigma_k = \sigma^2 I$

شکل بالا *training set* و همچنین خطوط دو توزیع گاوسی که با داده‌های هر یک از دو کلاس مطابقت دارند را نشان می‌دهد. توجه داشته باشید که دو توزیع گاوسی ماتریس کوواریانس Σ مشترک دارند ولی μ

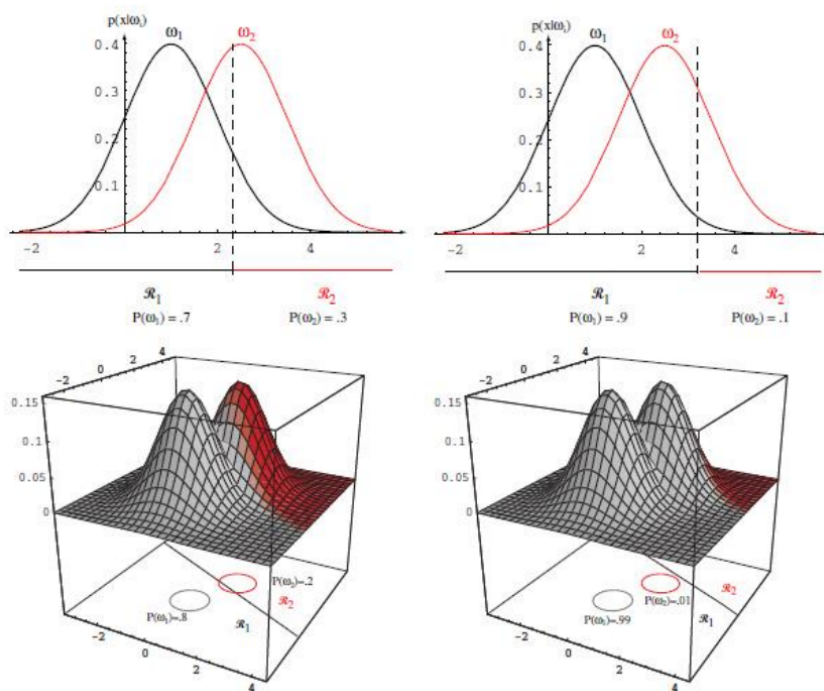
های متفاوت.

همچنین فرض داریم که $p(y = 1|x) = 0.5$



شکل ۱: حالت خاص Σ برابر

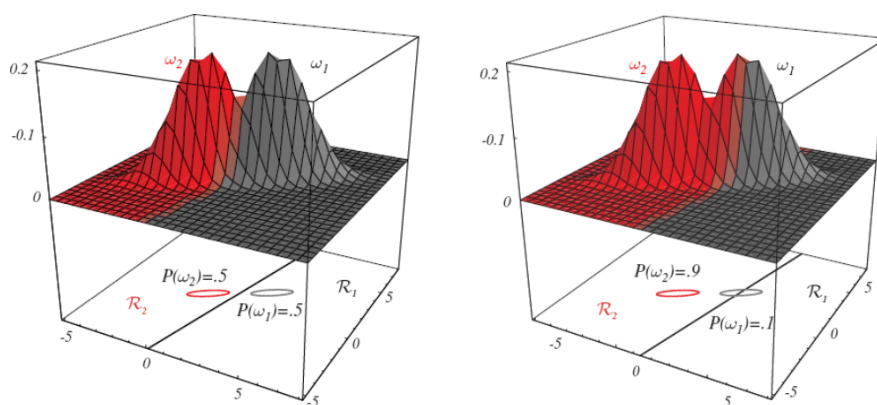
در این حالت خاص مرز تصمیم معمولاً در نقطه میانی بین دو مرکز کلاس قرار دارد، که منعکس کننده احتمال برابر برای هر کلاس است. این نشان دهنده نقطه ای است که هر دو کلاس به یک اندازه محتمل هستند و به عنوان جداکننده بین آنها عمل می کند و در نهایت یک تابع درجه ۱ است. به طور ساده میشود گفت که چون $prior$ ها مساوی اند $posterior$ برابر $likelihood$ است و دلیل وسط قرار گرفتن خط همین بود که در آن وسط $posterior$ ها هم همدیگر را قطع میکردند. حال فرض کنیم که $prior$ کلاس سیاه بیشتر هست در این صورت مرز به سمت قرمز میرود. در واقع تمایل به قرار گیری در کلاس سیاه بیشتر میشود به دلیل $prior$.



شکل ۲: حالت خاص Σ برابر

مورد خاص: $\Sigma_k = \Sigma$

در این حالت برای توزیع های گاوسی ای را داریم با ماتریس کواریانس های مساوی ولی دلخواه. ابرصفحه های تصمیم ما در این حالت خاص GDA معادل رویکرد LDA است. و اینکه ابرصفحه های تصمیم نباید بر خطی که میانه کلاس ها را به هم وصل می کند عمود باشند.



شکل ۳: حالت خاص $\Sigma_k = \Sigma$

Nave Bayes classifier

در GDA ، بردارهای ویژگی x بردارهای پیوسته و $real-valued$ بودند. بیایید اکنون در مورد یک الگوریتم یادگیری متفاوت صحبت کنیم که در آن x_j ها گسسته ارزیابی می شوند. برای مثال، ساخت فیلتر $spam$ ایمیل را با استفاده از ML در نظر بگیرید. در اینجا، می‌خواهیم پیام‌ها را بر اساس ایمیل‌های تجاری ناخواسته ($spam$) یا $non-spam$ طبقه‌بندی کنیم. پس از یادگیری انجام این کار، می‌توانیم ایمیل خوان خود را به‌طور خودکار پیام‌های $spam$ را فیلتر کرده و شاید آنها را در یک پوشه ایمیل جداگانه قرار دهیم. طبقه‌بندی ایمیل‌ها نمونه‌ای از مجموعه وسیع‌تری از مشکلات به نام طبقه‌بندی متن است.

فرض کنید یک $training set$ داریم (مجموعه‌ای از ایمیل‌ها که به عنوان $spam$ یا $non-spam$ برچسب گذاری شده‌اند). ما ساخت فیلتر $spam$ خود را با مشخص کردن ویژگی‌هایی که x_j برای نمایش ایمیل استفاده می‌شود، آغاز می‌کنیم. ما یک ایمیل را از طریق یک بردار ویژگی نشان خواهیم داد که طول آن برابر با تعداد کلمات در فرهنگ لغت است. به طور خاص، اگر یک ایمیل حاوی کلمه j -ام فرهنگ لغت باشد، $x_j = 1$ را تنظیم می‌کنیم. در غیر این صورت، اجازه می‌دهیم $x_j = 0$. برای مثال، بردار برای نشان

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{matrix} a \\ aardvark \\ aardwolf \\ \vdots \\ buy \\ \vdots \\ zygmurgy \end{matrix}$$

دادن ایمیلی استفاده می‌شود که حاوی کلمات " a " و " buy " است، اما نه " $aardwolf$ "، " $aardvark$ " یا " $zygmurgy$ ". مجموعه کلماتی که در بردار ویژگی کدگذاری شده‌اند، واژگان نامیده می‌شود. بنابراین بعد x برابر با اندازه واژگان است.

با انتخاب بردار ویژگی خود، اکنون می‌خواهیم یک مدل $generative$ بسازیم. بنابراین، باید $p(x|y)$ را مدل کنیم. اما اگر مثلاً واژگانی متشکل از ۵۰۰۰۰ کلمه داشته باشیم، آنگاه $x \in \{0, 1\}^{50000}$ یک بردار ۵۰۰۰۰ بعدی از ۰ و ۱ است، و اگر بخواهیم x را به طور صریح با توزیع $multinomial$ بر روی 2^{50000} نتیجه ممکن مدل کنیم، سپس ما به یک بردار پارامتر $(2^{50000} - 1)$ بعدی خواهیم رسید. این به وضوح تعداد پارامترهای زیادی است.

بنابراین برای مدل کردن $p(x|y)$ ، به یک فرض نیاز خواهیم داشت. یک فرضی وجود دارد به نام فرض *Nave Bayes* که میگوید با فرض کلاس، فیچر ها مستقل از هم تولید میشوند:

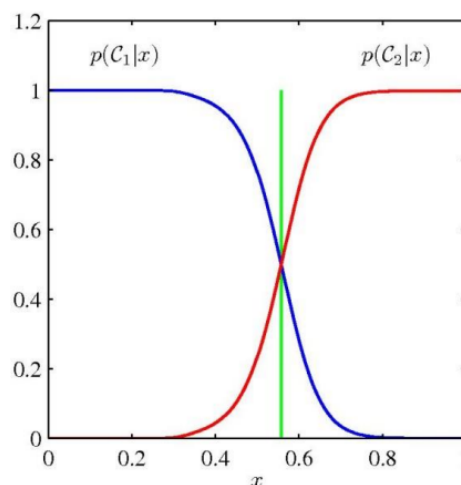
$$p(x|C_k) = p(x_1|C_k) \times p(x_2|C_k) \times \dots \times p(x_d|C_k)$$

برای مثال، اگر $y = 1$ به معنای ایمیل *spam* است. "buy" کلمه ۲۰۸۷ام و "price" کلمه ۳۹۸۳۱ام است. پس فرض می کنیم که اگر شما بگوییم $y = 1$ (که یک ایمیل خاص *spam* است)، دانش x_{2087} (دانش اینکه آیا "buy" در پیام ظاهر می شود) هیچ تاثیری بر باورهای شما در مورد ارزش x_{39831} ندارد. به طوری دیگر، این را می توان به صورت $p(x_{2087}|yx_{39831}) = p(x_{2087}|y)$ نوشت. (توجه داشته باشید که این همان چیزی نیست که بگوییم x_{2087} و x_{39831} مستقل هستند، که نوشته می شد $(=) p(x_{2087})$ بلکه فقط فرض می کنیم که x_{2087} و x_{39831} مستقل شرطی *given y* هستند) اکنون داریم:

$$\begin{aligned} p(x_1, \dots, x_{50000}|y) &= P(x_1|y)P(x_2|y, x_1)P(x_3|y, x_1, x_2) \dots p(x_{50000}|y, x_1, \dots, x_{49999}) \\ &= p(x_1|y)p(x_2|y)p(x_3|y) \dots p(x_{50000}|y) \\ &= \prod_{j=1}^{50000} p(x_j|y) \end{aligned}$$

Discriminative approach

در دیدگاه *Discriminative* مستقیم از توزیع *posterior* استفاده میکنیم برای هر کلاس C_k .



برای یک تعداد خوبی از توابع *conditional likelihood* توزیع *posterior* شبیه *logistic regression*

میشود. در واقع یکی از توابع مشهور برای مدل کردن *posterior* همین تابع *logistic regression* است.

یعنی فرضیمونو یک سیگموئید میگیریم از $w^t x$:

$$h(x; w) = \sigma(w^t x)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

داریم که:

$$x = \begin{bmatrix} 1 & x_1 & \dots & x_a \end{bmatrix}$$

$$w = \begin{bmatrix} w_0 & w_1 & \dots & w_a \end{bmatrix}$$

و x مساله در حالت تخمین میشه تخمین زدن مقادیر w ها. و x مقدار $h(x)$ بدیها بین ۰ و ۱ است و مقدار $y = 1$ را $given x$ نشان میدهد.