

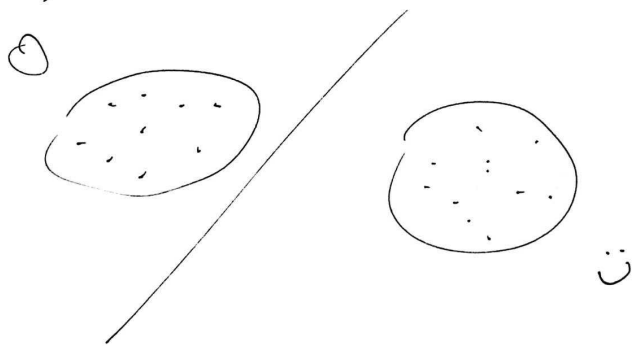
Generative Models

Sofar, Looked At

$$p(y|x; \theta)$$

For classification, really

i) A straight line to distinguish classes



Can do this differently

Model 0 \Rightarrow Can check X performance under Both Models

Model 1

Basically Learning $p(x|y)$ and $p(y)$ instead of $p(y|x)$

These are called generative (not discriminative) models

If y is class membership of 0 or 1 then

$$p(x|y=0)$$

have $P(y)$, $P(x|y)$, then we can use **Bayes' Theorem** to calculate $P(y|x)$.

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

$$P(x) = P(x|y=1)P(y=1) + P(x|y=0)P(y=0).$$

$$\text{Thus, } \arg\max_y P(y|x) = \arg\max_y \frac{P(x|y)P(y)}{P(x)} = \arg\max_y P(x|y)P(y)$$

Start with **Gaussian Discriminant Analysis (GDA)**; Need

Multivariate Normal Distributions.

In d dimensions, parameterized by mean vector $\mu \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$

with $\Sigma \geq 0$ symmetric and positive semidefinite

$(z^T \Sigma z \geq 0, \forall z \neq 0)_{z \in \mathbb{R}^d}$, then

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

↑
Determinant

Note for $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma)$
(vs μ, σ for 1D Normal),

$$E[\mathbf{X}] = \mu.$$

The covariance for a vector-valued RV \mathbf{X} is

$$\begin{aligned}\text{Cov}(\mathbf{X}) &= E[(\mathbf{X} - E[\mathbf{X}])(\mathbf{X} - E[\mathbf{X}])^T] \\ &= E[\mathbf{X}\mathbf{X}^T] - (E[\mathbf{X}](E[\mathbf{X}])^T),\end{aligned}$$

$$\text{for } \mathbf{X} \sim \mathcal{N}(\mu, \Sigma), \quad \text{Cov}(\mathbf{X}) = \Sigma$$

Gaussian Discriminant Analysis

Want to classify in situation where inputs are features \mathbf{X} , continuous random variable

Model $P(\mathbf{X}|\mathbf{y})$ as multivariate normal

i.e. $\mathbf{y} \sim \text{Bernoulli}(\theta)$

$$\left. \begin{aligned} \mathbf{X}|\mathbf{y}=0 &\sim \mathcal{N}(\mu_0, \Sigma) \\ \mathbf{X}|\mathbf{y}=1 &\sim \mathcal{N}(\mu_1, \Sigma) \end{aligned} \right\} \begin{array}{l} \text{same Covariance} \\ \text{Difference Mean} \end{array}$$

$$P(y) = \phi^y (1-\phi)^{1-y}$$

$$P(x|y=0) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0)\right)$$

$$P(x|y=1) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right)$$

And log likelihood:

$$\begin{aligned} \mathcal{L}(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n P(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\ &= \log \prod_{i=1}^n P(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) P(y^{(i)}; \phi) \end{aligned}$$

$$\phi = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\}$$

$$\mu_0 = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 0\}}$$

$$\mu_1 = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\}}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n \begin{pmatrix} x^{(i)} - \mu_{y^{(i)}} \end{pmatrix} \begin{pmatrix} x^{(i)} - \mu_{y^{(i)}} \end{pmatrix}^T$$

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^n (y^{(i)} \log \phi + (1 - y^{(i)}) \log(1 - \phi)) - \frac{1}{2} \log |\Sigma| + C$$

$$+ (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu)$$

Depends on $y^{(i)}$ if $\mu = 0$

$$\Rightarrow \mu = \mu_0 = 1 \Rightarrow \mu = \mu_1$$

Take partials wrt $\phi, \mu_0, \mu_1, \Sigma$; set = 0:

$$\frac{\partial \ell}{\partial \phi} = \sum_{i=1}^n \frac{y^{(i)}}{\phi} - \sum_{i=1}^n \frac{(1-y^{(i)})}{1-\phi} = 0$$

$$\Rightarrow (1-\phi) \sum_{i=1}^n y^{(i)} = \phi \sum_{i=1}^n (1-y^{(i)})$$

$$\sum_{i=1}^n y^{(i)} - \phi \sum_{i=1}^n y^{(i)} = n\phi - \phi \sum_{i=1}^n y^{(i)}$$

$$\phi = \frac{1}{n} \sum_{i=1}^n y^{(i)} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}$$

Note $\nabla_A \text{tr} A B A^T C = C A B + C^T A B^T$, Σ^{-1} is symmetric

$$\nabla_{\mu_0} \ell = \nabla_{\mu_0} \sum_{i=1}^n \underbrace{\text{tr}(\mu^{(i)} - \mu)^T \Sigma^{-1} (\mu^{(i)} - \mu)}_{1 \times 1} \mathbb{I}$$

$$= - \sum_{i=1}^n \left[\mathbb{I}(\mu^{(i)} - \mu_0)^T \Sigma^{-1} + \mathbb{I}(\mu^{(i)} - \mu_0) (\Sigma^{-1})^T \right] \mathbb{1}\{y^{(i)}=0\} = 0$$

$$\Rightarrow -2 \Sigma^{-1} \underbrace{\left[\sum_{i=1}^n (\mu^{(i)} - \mu_0) \mathbb{1}\{y^{(i)}=0\} \right]}_{=0} = 0 \quad \Sigma^{-1} \neq 0 \text{ is invertible}$$

$$\Rightarrow \sum_{i=1}^n (x^{(i)} - \mu_0) \mathbb{1}_{\{y^{(i)} = 0\}} = 0$$

$$\Rightarrow \mu_0 = \frac{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 0\}} x^{(i)}}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)} = 0\}}}$$

Now, find Σ : Note $\frac{1}{|\Sigma|} = |\Sigma^{-1}|$

Instead of $x^{(i)}$, write x , μ instead of μ_0 or μ_1 .

$$\text{tr } ABC = \text{tr } CAB$$

Take the derivative (Dropping Terms That Don't Matter)

$$\mathcal{L} = \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n \text{tr} (x - \mu)(x - \mu)^T \Sigma^{-1}$$

$$l(\varphi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^n \left[y^{(i)} \log \varphi + (1 - y^{(i)}) \log (1 - \varphi) - \frac{1}{2} \log |\Sigma| + C + (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) \right]$$

Ignoring terms that go to 0 when we take gradient,

$$l = \frac{n}{2} \log |\Sigma^{-1}| - \frac{1}{2} \sum_{i=1}^n \text{tr} (x - \mu_{y^{(i)}})(x - \mu_{y^{(i)}})^T \Sigma^{-1}$$

$$\text{Note } \nabla_A \log |A| = (A^{-1})^T, \quad \nabla_B \text{tr} AB = A^T.$$

Let $\Lambda = \Sigma^{-1}$, then

$$\nabla_{\Lambda} l = \nabla_{\Lambda} \left[\frac{n}{2} \log |\Lambda| - \frac{1}{2} \sum_{i=1}^n \text{tr} (x - \mu)(x - \mu)^T \Lambda \right]$$

$$= \frac{n}{2} \Lambda^{-1} - \frac{1}{2} \sum_{i=1}^n (x - \mu)(x - \mu)^T \stackrel{\text{set}}{=} 0 \quad \leftarrow \text{for critical point}$$

$$n \Lambda^{-1} = \sum_{i=1}^n (x - \mu)(x - \mu)^T$$

$$n \Sigma = \sum_{i=1}^n (x - \mu)(x - \mu)^T$$

$$\Rightarrow \Sigma = \frac{1}{n} \sum_{i=1}^n (x - \mu_{y^{(i)}})(x - \mu_{y^{(i)}})^T$$

Look at the following as function of x :

$$P(x=1 | x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-\Theta^T x)}$$



This looks like logistic regression!

Thus, why would I use **GDA** vs. **Log. Reg.** in any particular setting?

Give different outputs

Turns out, $P(x|y)$ multivariate Gaussian w/ $\Sigma \Rightarrow$

$P(y|x)$ is logistic.

But $P(y|x)$ logistic $\nRightarrow P(x|y)$ mult. Gaussian,

$x|y=0 \sim \text{Poisson}(\lambda_0) \Rightarrow P(y|x)$ Logistic

$x|y=1 \sim \text{Poisson}(\lambda_1)$

If your data looks \approx Gaussian,

then GDA will work better.

Logistic Reg. cares less about underlying distributions, so is more robust.

Naive Bayes

Recall: Sum Rule: $P(x) = \sum_y P(x, y)$

Product Rule: $P(x, y) = P(y|x)P(x)$

Want a model that assumes x is discrete and y is discrete classification.

Going to introduce the "Naive Bayes" model for binary classification.

Sample problem - SPAM filtering

Think of an email x as a vector of whether or not a word appears in an email:

$$x = \left[\begin{array}{c} 0 \\ 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ 0 \end{array} \right] \left. \begin{array}{l} \text{win} \\ \text{pure} \\ \text{gold} \\ \text{xxx} \\ \text{office} \\ \text{hours} \\ \text{copier} \\ \vdots \\ \text{toast} \end{array} \right\} \begin{array}{l} \text{words called} \\ \text{vocabulary} \end{array} \in \mathbb{R}^d$$

d is # of words

Goal is to model $P(x|y)$, but if d is large size of x space goes like is 2^d .

Model this w/multinomial need $2^d - 1$ parameters. Thus, assume, to reduce size of model, that x_i 's are conditionally independent given y (Not Ind.). This means if I already know what class an email is in, then appearance of individual words is independent.

i.e. if I know it is spam, knowing "pure" occurs does NOT tell me anything about the occurrence of "gold".

$$\text{Mathematically, } P(x_7|y) = P(x_7|y, x_3)$$

$$\stackrel{\text{Not}}{\Rightarrow} P(x_7) = P(x_7|x_3)$$

How is this helpful?

$$P(x_1, \dots, x_d | y) = \underset{\substack{\uparrow \\ \text{Always true}}}{P(x_1 | y)} P(x_2 | y, x_1) \cdot \dots \cdot P(x_d | y, x_1, \dots, x_{d-1})$$

Naive Bayes' Assumption

$$= P(x_1 | y) P(x_2 | y) \cdot \dots \cdot P(x_d | y)$$

$$= \prod_{j=1}^d P(x_j | y)$$

What are the learnable parameters?

$$\varphi_{j|y=1} = P(x_j = 1 | y = 1), \quad \varphi_y = P(y = 1)$$

$$\varphi_{j|y=0} = P(x_j = 1 | y = 0)$$

For training set $\{(x^{(i)}, y^{(i)})\}, i = 1, \dots, n\}$

$$\mathcal{L}(\varphi_y, \varphi_{j|y=0}, \varphi_{j|y=1}) = \prod_{i=1}^n P(x^{(i)}, y^{(i)})$$

Maximized w/ ..

Maximized with

$$\phi_j | y=1 = \frac{\sum_{i=1}^n \mathbb{1}\{x_j^{(i)}=1 \wedge y^{(i)}=1\}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}}$$

$$\phi_j | y=0 = \frac{\sum_{i=1}^n \mathbb{1}\{x_j^{(i)}=1 \wedge y^{(i)}=0\}}{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\}}$$

$$\phi_y = \frac{\sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}}{n}$$

Given a new email, how do I classify it?

Calculate $P(y=1|x)$, $P(y=0|x)$.

choose Bigger One

$$P(y=1|x) = \frac{P(x|y=1)P(y=1)}{P(x)}$$

$$= \frac{\left(\prod_{i=1}^d P(x_i|y=1) P(y=1) \right)}{\left(\prod_{j=1}^d P(x_j|y=1) \right) P(y=1) + \left(\prod_{j=1}^d P(x_j|y=0) \right) P(y=0)}$$

What if Naive Bayes never saw a word in the training data "pure"

$$\varphi_j | y=0 = \varphi_{j=y=1} = 0$$

Now this word appears in an email you want to classify $\Rightarrow p(y=1 | x) = \frac{0}{0}$. What to do?

For estimating a multinomial distribution $\varphi_j = p(z=j)$ with n observations and k classes

$\{z^{(1)}, \dots, z^{(n)}\}$, max likelihood is

$$\varphi_j = \frac{\sum_{i=1}^n \mathbb{1}\{z^{(i)} = j\}}{n} \quad \text{to not allow } \varphi_j = 0.$$

Use Laplace Smoothing:

$$\varphi_j = \frac{1 + \sum_{i=1}^n \mathbb{1}\{x^{(i)} = j\}}{k + n} \quad (\text{Note } \sum \varphi_j = 1)$$

Thus, applying this for Naive Bayes:

$$\varphi_{j|y=1} = \frac{1 + \sum_{i=1}^n \mathbb{1}\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{2 + \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\}}$$

$$\varphi_{j|y=0} = \frac{1 + \sum_{i=1}^n \mathbb{1}\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{2 + \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 0\}}$$