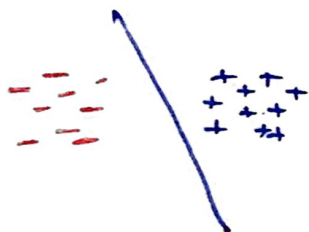


Objectives:

- Understand what Generative Models are and how they work
- Understand estimation and prediction phases of generative models
- Derive a relation connecting generative and discriminative models
- Derive Maximum Likelihood Estimates (MLE) for multinomial and Gaussian generative models.

Generative vs Discriminative models:

Whenever we talked about classification, the pictures that you had in mind is that your classifier had a training data, let's say positives and negative instances. The job of the discriminative model was to find a separator that discriminates between these points.



The approach that we will be taking is actually trying to understand what is the structure of these \oplus and \ominus classes.

If I can understand in some probabilistic term, the structure of positive and negative examples maybe I can do discrimination in a different way that we've done when we're primarily looking at the separator.

Two questions to ask about any generative model:

1. how we estimate this model? Estimation question
2. how we actually do prediction? Prediction question

Given our training data, we will fit probability distribution for the negative class and for the positive class. And by comparing these probabilities, we can actually induce what is the right label.

Simple Multinomial Generative model:

We will talk about multinomial and in your head, you can think about text documents and what our models will do, they will generate documents. Some parameters:

$$p(w|\theta) = \theta_w$$

\uparrow a word \uparrow parameters of a model

captures what's the likelihood of me selecting, generating certain words given all the possibilities.

Selecting words independently from each other.

$$\theta_w \geq 0, \sum_{w \in W} \theta_w = 1$$

likelihood function:

let's say we have θ_w , how do I compute the likelihood of generating a document?

$$P(D|\theta) = \prod_{i=1}^n \theta_{w_i}$$

we have our document
probability of every word in our document

$$P(D|\theta) = \prod_{w \in W} \theta_w^{\text{count}(w)}$$

the words in our vocabulary

let's take an example: we have two words in our vocabulary

$$W = \{\text{cat}, \text{dog}\}$$

model 1 which takes θ , $\theta_{\text{cat}} = 0.3$, $\theta_{\text{dog}} = 0.7$

model 2 which takes θ' , $\theta'_{\text{cat}} = 0.9$, $\theta'_{\text{dog}} = 0.1$

$D = \{\text{cat}, \text{cat}, \text{dog}\}$ we can compute the likelihood of these documents generated by 1 and 2 model

$$P(D|\theta) = (0.3)^2 \times (0.7)$$

$$P(D|\theta') = (0.9)^2 \times (0.1)$$

How can we utilize our training data to find the best parameters?

we will use Maximum likelihood, and will make assumption, that the best parameters are the parameters which give the highest likelihood to our data.

Find θ 's which maximize $P(D|\theta) = \max_{\theta} \prod_{w \in W} \theta_w^{\text{count}(w)}$

$$\log \prod_{w \in W} \theta_w^{\text{count}(w)} = \sum_{w \in W} \text{count}(w) \log \theta_w$$

$$W = \left\{ \begin{array}{cc} 0 & 1 \\ \text{cat} & \text{dog} \end{array} \right\} \quad \left| \quad \frac{\partial}{\partial \theta} \text{count}(0) \cdot \log(\theta) + \text{count}(1) \cdot \log(1-\theta) \right.$$

$$\theta_{\text{cat}} = \theta_0 = \theta$$

$$\theta_{\text{dog}} = \theta_1 = 1 - \theta$$

$$\Rightarrow \frac{\text{count}(0)}{\theta} - \frac{\text{count}(1)}{1-\theta} = 0$$

$$(1-\theta) \text{count}(0) - \theta \text{count}(1) = 0$$

$$\hat{\theta} = \frac{\text{count}(0)}{\text{count}(1) + \text{count}(0)}$$

Prediction:

60

$$\text{MLE for Multinomial Distribution: } \hat{\theta}_w = \frac{\text{count}(w)}{\sum_{w' \in W} \text{count}(w')}$$

Using the estimation techniques described earlier, we can find θ^+ and θ^- that they will give the highest likelihood to the points.

The question is, if I give you a new document, how do you know to which class it belongs? I could look at the likelihood that this document was generated by plus side and minus side.

$$\log \frac{P(D|\theta^+)}{P(D|\theta^-)} \quad \text{For now let's assume that the prior likelihood of the } \ominus \text{ and } \oplus \text{ class are exactly the same.}$$

$$\text{In this case: } \log \frac{P(D|\theta^+)}{P(D|\theta^-)} = \begin{cases} \geq 0 \rightarrow + \\ < 0 \rightarrow - \end{cases}$$

\rightarrow class conditional distribution.

$$= \log P(D|\theta^+) - \log P(D|\theta^-)$$

$$= \log \prod_{w \in W} \theta_w^{+\text{count}(w)} - \log \prod_{w \in W} \theta_w^{-\text{count}(w)}$$

$$= \sum_{w \in W} \text{count}(w) \log \theta_w^+ - \sum_{w \in W} \text{count}(w) \log \theta_w^-$$

$$= \sum_{w \in W} \text{count}(w) \log \frac{\theta_w^+}{\theta_w^-} \quad \tilde{\theta}_w$$

$$= \sum_{w \in W} \text{count}(w) \cdot \tilde{\theta}_w$$

Prior, Posterior and Likelihood:

Sometimes we might have some prior knowledge and we want to take advantage of it.

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad \text{Bayesian Rule}$$

$P(B)$ is the Prior

$$\text{posterior} \rightarrow \underbrace{P(y=+|D)} = \frac{P(D|\theta^+) \cdot \overbrace{P(y=+)}^{\text{Prior}}}{P(D)} \quad \text{What is the likelihood that I'm going to assign to document } D \text{ it's label } y \text{ which is } \oplus?$$

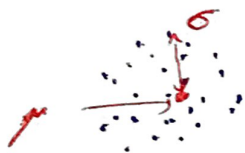
$$\log \frac{P(y=+|D)}{P(y=-|D)} = \log \frac{P(D|\theta^+) \cdot P(y=+)}{P(D|\theta^-) \cdot P(y=-)}$$

$$\log \left[\frac{P(D|\theta^+)}{P(D|\theta^-)} \right] \left[\frac{P(y=+)}{P(y=-)} \right] = \log \frac{P(D|\theta^+)}{P(D|\theta^-)} + \log \frac{P(y=+)}{P(y=-)} \quad \leftarrow \tilde{\theta}_0$$

$$= \sum_{w \in W} \text{count}(w) \cdot \tilde{\theta}_w + \tilde{\theta}_0 \quad \text{We translate it to a linear classifier}$$

Gaussian Generative models:

Now we are going to vectors in \mathbb{R}^d and find distribution that is a very natural fit for describing this data. $X \in \mathbb{R}^d$



We will describe these points by two parameters:

1. Where is the center of the cloud? μ
2. How dispersed is this cloud? σ^2

$$P(X|y, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} \exp\left(-\frac{1}{2\sigma^2} \|X - \mu\|^2\right)$$

Multivariate Gaussian Random vector:

A random vector $X = (X^{(1)}, \dots, X^{(d)})^T$ is a **Gaussian vector**, or **multivariate Gaussian** or **normal variable**, if any linear combination of its components is a (univariate) Gaussian variable or a constant (a "Gaussian" variable with zero variance) i.e., if $\alpha^T X$ is (univariate) Gaussian or constant for any constant non-zero vector $\alpha \in \mathbb{R}^d$

The distribution of X , the d -dimensional **Gaussian** or **normal distribution**, is completely specified by the vector mean $\mu = E[X] = (E[X^{(1)}], \dots, E[X^{(d)}])^T$ and the $d \times d$ covariance matrix Σ . If Σ is invertible, then the pdf of X is:

$$f_X(X) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(X-\mu)^T \Sigma^{-1}(X-\mu)} \quad X \in \mathbb{R}^d$$

Where $\det(\Sigma)$ is the determinant of the Σ , which is positive when Σ is invertible

if $\mu = 0$ and Σ is the identity matrix, then X is called a **Standard normal random vector**

MLE for Gaussian Distribution:

62

In this problem we will derive the maximum likelihood estimator for a Gaussian model.

$$f_X(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\|x-\mu\|^2}{2\sigma^2}}$$

Let $S_n = \{x^{(1)}, x^{(2)}, \dots, x^{(n)}\}$ be i.i.d random variables following a Gaussian distribution with mean μ and variance σ^2 , then Their joint probability density function is given by:

$$\prod_{t=1}^n P(x^{(t)}|\mu, \sigma^2) = \prod_{t=1}^n \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\|x^{(t)}-\mu\|^2}{2\sigma^2}}$$

Taking logarithm of the above function, we get:

$$\begin{aligned} \log P(S_n|\mu, \sigma^2) &= \log \left(\prod_{t=1}^n \frac{1}{(2\pi\sigma^2)^{d/2}} e^{-\frac{\|x^{(t)}-\mu\|^2}{2\sigma^2}} \right) \\ &= \sum_{t=1}^n \log \frac{1}{(2\pi\sigma^2)^{d/2}} + \sum_{t=1}^n \log e^{-\frac{\|x^{(t)}-\mu\|^2}{2\sigma^2}} \\ &= \sum_{t=1}^n -\frac{d}{2} \log(2\pi\sigma^2) + \sum_{t=1}^n \log e^{-\frac{\|x^{(t)}-\mu\|^2}{2\sigma^2}} \\ &= -\frac{nd}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{t=1}^n \|x^{(t)}-\mu\|^2 \end{aligned}$$

$$\frac{\partial \log P(S_n|\mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{t=1}^n (x^{(t)} - \mu) = 0$$

$$\hat{\mu} = \frac{\sum_{t=1}^n x^{(t)}}{n}$$

$$\frac{\partial \log P(S_n|\mu, \sigma^2)}{\partial \sigma^2} = -\frac{nd}{2\sigma^2} + \frac{\sum_{t=1}^n \|x^{(t)}-\mu\|^2}{2(\sigma^2)^2}$$

$$\hat{\sigma}^2 = \frac{\sum_{t=1}^n \|x^{(t)}-\mu\|^2}{nd}$$