Bayes classifier, Decision boundary

Aarti Singh

Machine Learning 10-315 Sept 1, 2021



In-person Expectations

- Please do not come to class if you feel sick or have symptoms of a potentially contagious illness.
- This holds <u>even if you know you do not have COVID-19</u> but have symptoms that may be a sign of other contagious illnesses, such as a <u>cold or flu</u>.
- Coming to class when sick is NOT a sign of hard work. Be responsible, care for your classmates and campus community.

Announcements

- Canvas fixed
- Late days total 4, no more than 1 for a QnA, no more than 2 for a HW
- Tentative HW due dates have been posted
- QnA1 out today
 - due in 1 week, Sept 8 11:59 pm ET

•	Office hours	Day	Time	Location	Staff
		Mondays	1:30 pm	TBA	Christina
		Tuesdays	3:30 pm	TBA	Spoorthi
		Wednesdays	2:30 pm	TBA	Haitian
		Thursdays	10:30 am	Zoom (link on Canvas)	Aarti

Notion of "Features aka Attributes"





remember to wake up when class ends

wake ends to class remember up when

How to represent inputs mathematically?

- Document vector X > Ideas?
 - list of words (different length for each document)
 - frequency of words (length of each document = size of vocabulary), also known as Bag-of-words approach

Misses out context!!

list of n-grams (n-tuples of words)

Why might this be limited?

Notion of "Features aka Attributes"

Input $X \in \mathcal{X}$



Input $X \in \mathcal{X}$



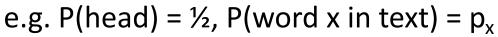
How to represent inputs mathematically?

- Image X = intensity/value at each pixel, fourier transform values, SIFT etc.
- Market information X = daily/monthly? price of share for past
 10 years

Distribution of Inputs

Input
$$X \in \mathcal{X}$$

Discrete Probability Distribution P(X) = P(X=x)





Probabilities in a distribution sum to 1

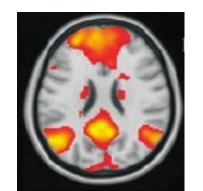
$$\sum_{x} P(X=x) = 1$$
 $P(tail) = 1 - p(head), \sum_{x} p_{x} = 1$

Continuous Probability density p(x)

Probability density integrate to 1

$$\int p(x)dx = 1$$

$$P(a \le X \le b) = \int_{a}^{b} p(x) dx$$

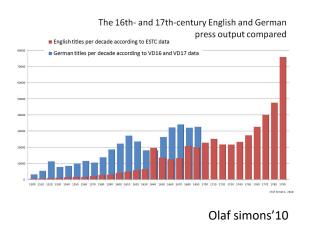


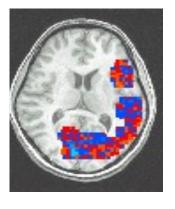
Distributions in Supervised tasks

Input
$$X \in \mathcal{X}$$

 Distribution learning also arises in supervised learning tasks e.g. classification

P(X = x | Y = y) Distribution of words in 'news' documents Distribution of brain activity under 'stress'

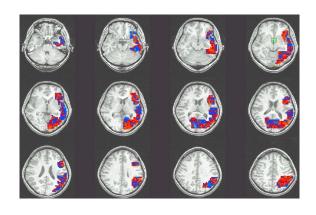




P(Y = y | X = x) Distribution of topics given document

Classification

Goal: Construct **prediction rule** $f: \mathcal{X} \to \mathcal{Y}$





High Stress
Moderate Stress
Low Stress

Input feature vector, X

Label, Y

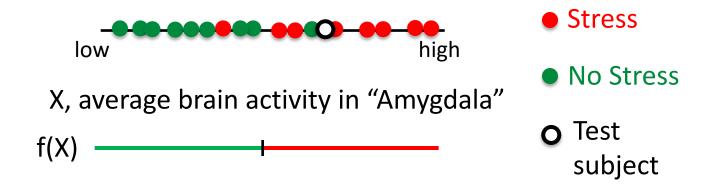
In general: label Y can belong to more than two classes

X is multi-dimensional (many features represent an input)

But lets start with a simple case:

label Y is binary (either "Stress" or "No Stress") X is average brain activity in the "Amygdala"

Binary Classification



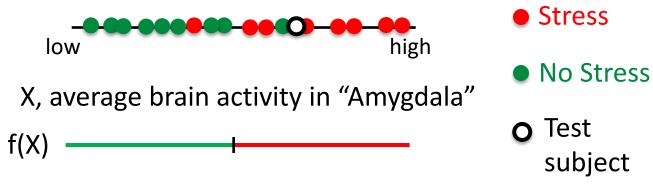
Model X and Y as random variables with joint distribution P_{XY}

Training data $\{X_i, Y_i\}_{i=1}^n \sim iid (independent)$ and identically distributed) samples from P_{XY}

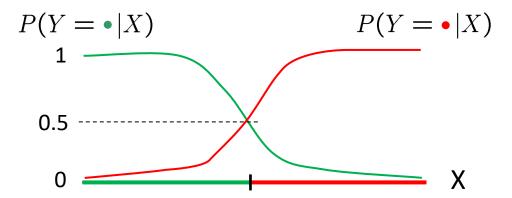
Test data $\{X,Y\}$ ~ iid sample from P_{XY}

Training and test data are independent draws from **same** distribution

Bayes Classifier



Model X and Y as random variables



For a given X, f(X) = label Y which is more likely

$$f(X) = \arg \max_{Y=y} P(Y=y|X=x)$$

Optimality of Bayes Classifier

Bayes Rule

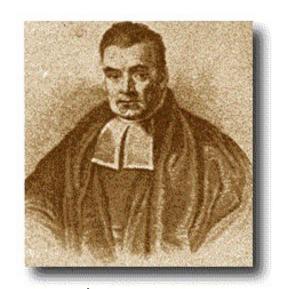
Bayes Rule:
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

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

To see this, recall:

$$P(X,Y) = P(X|Y) P(Y)$$

$$P(Y,X) = P(Y|X) P(X)$$



Thomas Bayes

Bayes Classifier – equivalent form

Bayes Rule:
$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

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

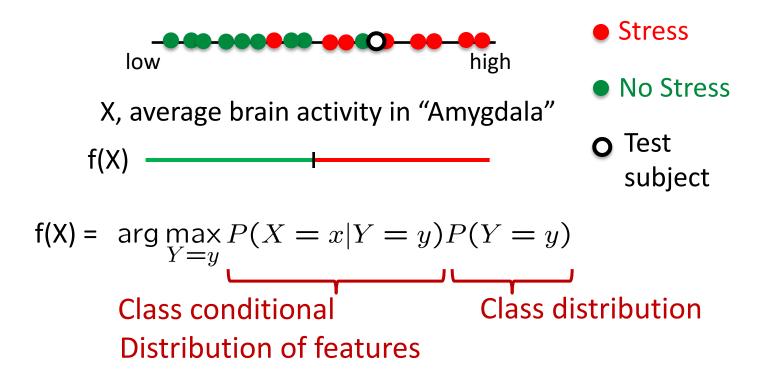
Bayes classifier:

$$f(X) = \arg \max_{Y=y} P(Y=y|X=x)$$

$$= \arg \max_{Y=y} P(X=x|Y=y)P(Y=y)$$

Class conditional Distribution of class Distribution of features

Bayes Classifier

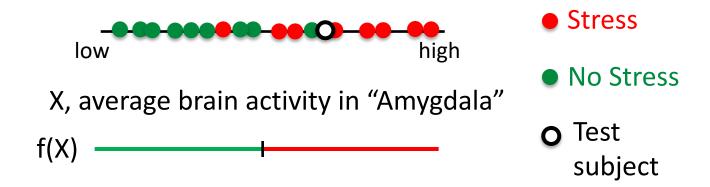


We can now consider appropriate distribution models for the two terms:

Class distribution P(Y=y)

Class conditional distribution of features P(X=x|Y=y)

Modeling class distribution



Modeling Class distribution $P(Y=y) = Bernoulli(\theta)$

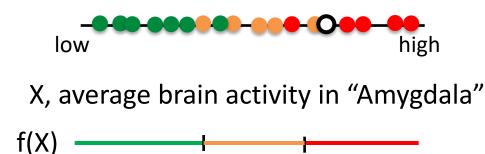
$$P(Y = \bullet) = \theta$$

$$P(Y = 0) = 1 - \theta$$

Like a coin flip



Modeling class distribution



- High Stress
- Moderate Stress
- Low Stress
- o Test subject
- ➤ How do we model multiple (>2) classes?

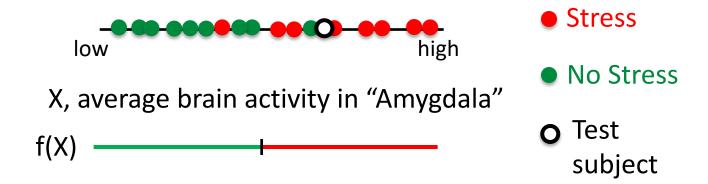
Modeling Class distribution $P(Y) = Multinomial(p_H, p_M, p_L)$

Like a dice roll



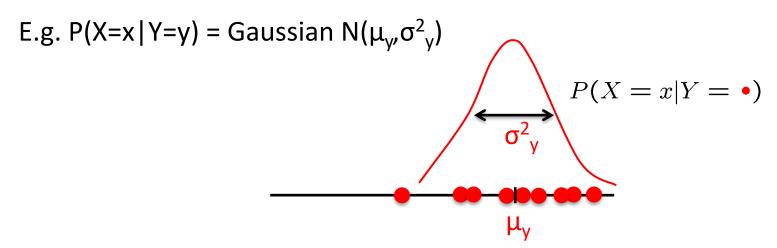
$$p_{H} + p_{M} + p_{I} = 1$$

Modeling class conditional distribution of features



Modeling class conditional distribution of feature P(X=x|Y=y)

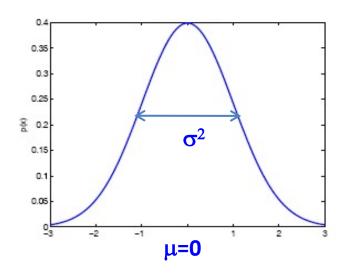
What distribution would you use?

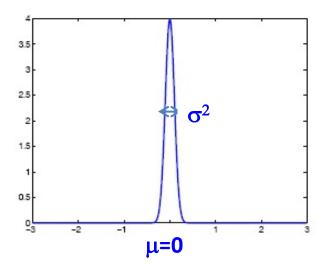


1-dim Gaussian distribution

X is Gaussian $N(\mu,\sigma^2)$

$$P(X = x | \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$





Why Gaussian?

- Properties
 - Fully Specified by first and second order statistics
 - Uncorrelated ⇔ Independence
 - X, Y Gaussian => aX+bY Gaussian
 - <u>Central limit theorem:</u> if X_1 , ..., X_n are any iid random variables with mean μ and variance $\sigma^2 < \infty$ then

$$\sqrt{n}\left(\frac{1}{n}\sum_{i=1}^{n}X_{i}-\mu\right)\sim N(0,\sigma^{2})$$

1-dim Gaussian Bayes classifier

$$f(X) = \arg \max_{Y=y} P(X = x | Y = y) P(Y = y)$$

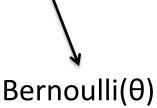
Learn parameters θ , μ_y , σ_v from data

Class conditional

Distribution of features

Gaussian(μ_y , σ^2_y)

Class distribution



$$P(Y = \bullet)P(X = x|Y = \bullet)$$

$$P(Y = \bullet)P(X = x|Y = \bullet)$$

1-dim Gaussian Bayes classifier

f(X) =
$$\underset{Y=y}{\operatorname{arg\,max}} P(X = x | Y = y) P(Y = y)$$

Class conditional Class distribution

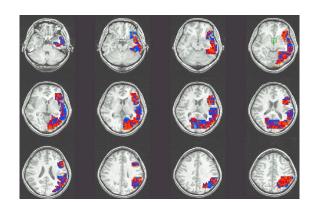
Distribution of features

Gaussian($\mu_{\text{W}}, \sigma^2_{\text{V}}$) Bernoulli(θ)

What decision boundaries can we get in 1-dim?

 $P(Y = \bullet)P(X = x|Y = \bullet)$ $P(Y = \bullet)P(X = x|Y = \bullet)$ 21

d-dimensional inputs





High Stress
Moderate Stress
Low Stress

Input feature vector, X

Label, Y

Modeling class conditional distribution of feature P(X=x|Y=y)

What distribution would you use?

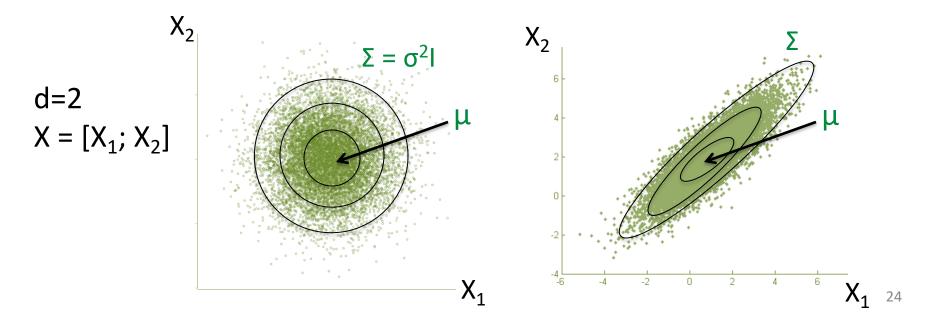
E.g. $P(X=x|Y=y) = Gaussian N(\mu_{v}, \Sigma_{v})$

d-dim Gaussian distribution

X is Gaussian $N(\mu, \Sigma)$

 μ is d-dim vector, Σ is dxd dim matrix

$$P(X = x | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right),$$



d-dim Gaussian Bayes classifier

$$f(X) = \arg \max_{Y=y} P(X = x | Y = y) P(Y = y)$$

Learn parameters θ , μ_y , Σ_v from data

Class conditional Distribution of features

Gaussian(μ_{v} , Σ_{v})

Class distribution

¥ Bernoulli(θ)

$$P(Y = \bullet)P(X = x|Y = \bullet)$$

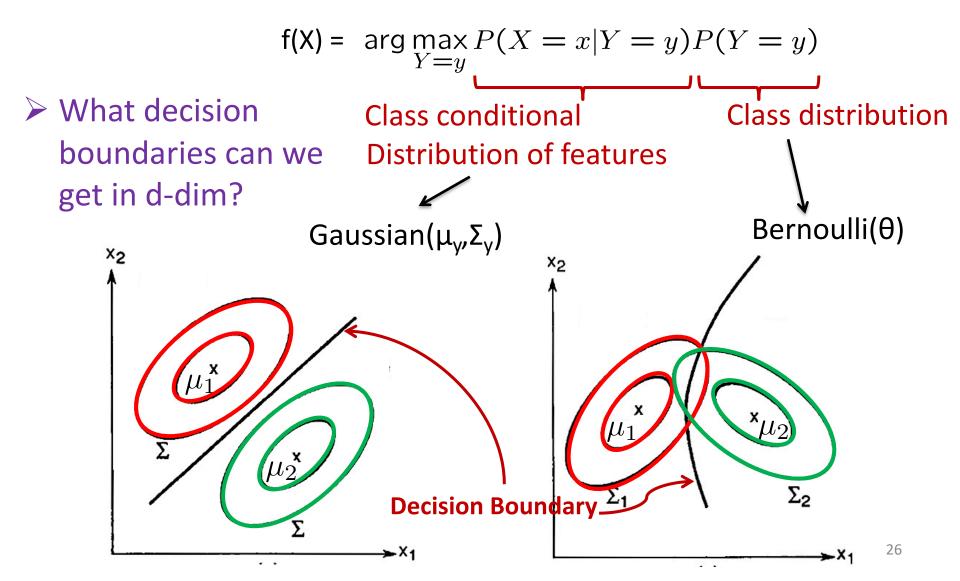
$$\mu_1^{\mathsf{x}}$$

$$\mu_2^{\mathsf{x}}$$

$$P(Y = \bullet)P(X = x|Y = \bullet)$$

$$\Sigma$$

d-dim Gaussian Bayes classifier



Decision Boundary of Gaussian Bayes

Decision boundary is set of points x: P(Y=1|X=x) = P(Y=0|X=x)

Compute the ratio

$$1 = \frac{P(Y=1|X=x)}{P(Y=0|X=x)} = \frac{P(X=x|Y=1)P(Y=1)}{P(X=x|Y=0)P(Y=0)}$$

$$= \sqrt{\frac{|\Sigma_0|}{|\Sigma_1|}} \exp\left(-\frac{(x-\mu_1)^{\mathsf{T}}\Sigma_1^{-1}(x-\mu_1)}{2} + \frac{(x-\mu_0)^{\mathsf{T}}\Sigma_0^{-1}(x-\mu_0)}{2}\right) \frac{\theta}{1-\theta}$$

In general, this implies a quadratic equation in x. But if $\Sigma_1 = \Sigma_0$, then quadratic part cancels out and decision boundary is linear.

Recap

Bayes classifier – assumes P_{XY} known, optimal for 0/1 loss

$$f(X) = \arg\max_{Y=y} P(Y=y|X=x)$$

$$= \arg\max_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \limsup_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \max_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \max_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \max_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \max_{Y=y} P(X=x|Y=y)P(Y=y)$$

$$= \min_{Y=y} P(X=x|Y=y)P(X=y)$$

$$= \min_{Y=y} P(X=x|Y=y)$$

- Gaussian Bayes classifier assumes
 Class distribution is Bernoulli/Multinomial
 Class conditional distribution of features is Gaussian
- Decision boundary (binary classification)

Glossary of Machine Learning

- Feature/Attribute
- iid
- Bayes classifier
- Class distribution
- Class conditional distribution of features
- Decision boundary