

Decision boundary $\{x: P(Y=1|X=x) = P(Y=0|X=x)\}$ ←

d-dim Gaussian Bayes classifier

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$f(x) = \arg \max_{Y=y} \underbrace{P(X=x|Y=y)}_{\text{Class conditional}} \underbrace{P(Y=y)}_{\text{Class distribution}}$$

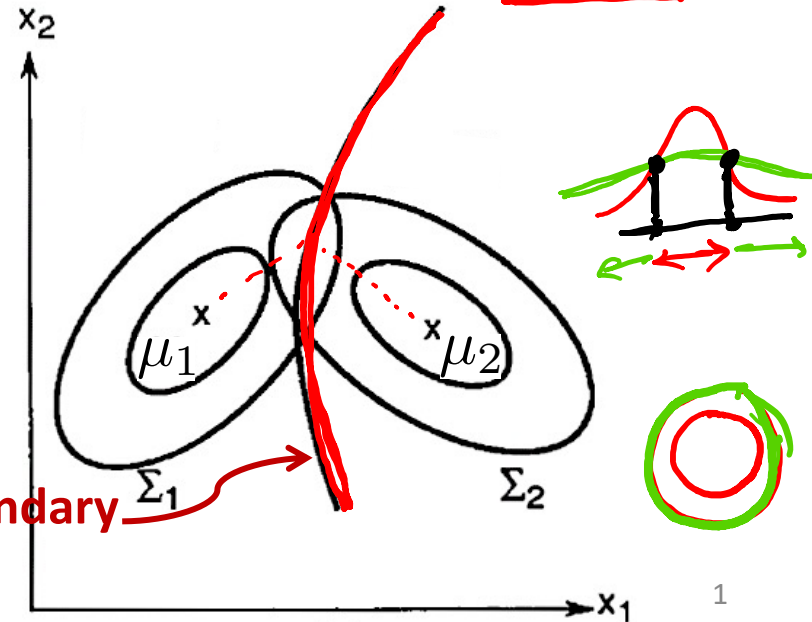
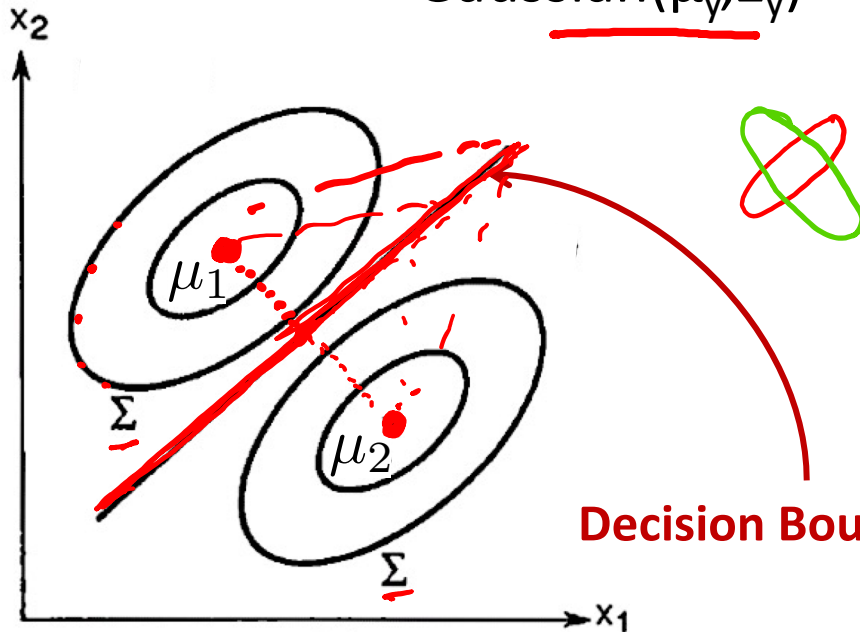
➤ What decision boundaries can we get in d-dim?

Class conditional
Distribution of features

Class distribution

Gaussian(μ_y, Σ_y)

Bernoulli(θ)



Decision Boundary

Decision Boundary of Gaussian Bayes

- Decision boundary is set of points x : $P(Y=1|X=x) = P(Y=0|X=x)$
- By Bayes theorem, equivalent to x :

$$\left\{ x: p(X=x|Y=1) P(Y=1) = p(X=x|Y=0) P(Y=0) \right\}$$

Lets find the decision boundary.

If class distribution is $P(Y=1) = \text{Ber}(\theta)$ and class conditional feature distribution $P(X=x|Y=y)$ is d -dim Gaussian $N(\mu_y, \Sigma_y)$

$$P(X = x|Y = y) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_y|}} \exp\left(-\frac{(x - \mu_y)^\top \Sigma_y^{-1} (x - \mu_y)}{2}\right)$$

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

$$-x^T \Sigma_1^{-1} x$$
$$+ x^T \Sigma_0^{-1} x$$

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

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

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.

Glossary of Machine Learning

- Feature/Attribute
- iid
- Bayes classifier
- Class distribution
- Class conditional distribution of features
- Estimator – hat notation
- MLE
- Decision boundary

Some notes

- Recitation Friday Sept 18
 - Recap of MLE/MAP hands on
 - Naïve Bayes application
 - Linear algebra and multi-variate calculus
- HW1 due date -> Sept 25

Naïve Bayes

Learning Distributions (MAP)

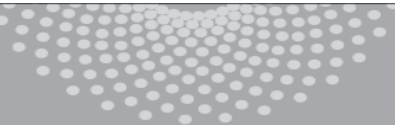
Aarti Singh

Machine Learning 10-315

Sept 16, 2020

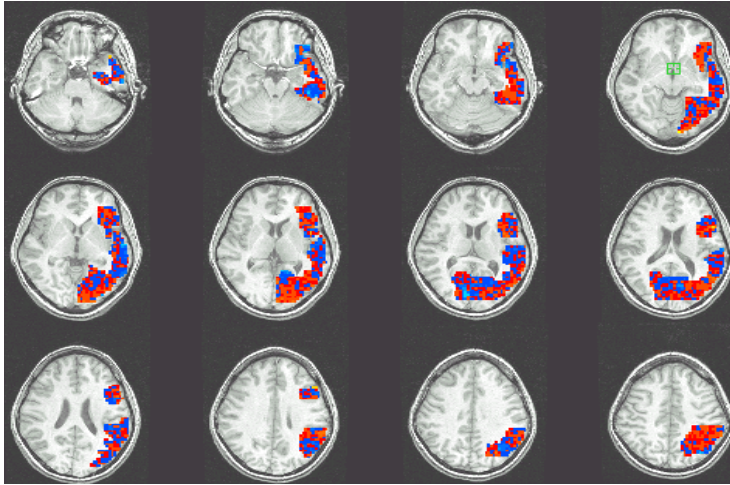


MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Multi-class, multi-dimensional classification – Continuous features



Input feature vector, X



High Stress
Moderate Stress
Low Stress

Label, Y

We started with a simple case:

label Y is binary (either “Stress” or “No Stress”)

X is average brain activity in the “Amygdala”

In general: label Y can belong to $K > 2$ classes

X is multi-dimensional $d > 1$ (average activity in all brain regions)

How many parameters do we need to learn (continuous features)?

Class probability: $K=3$

$P(Y = y) = p_y$ for all y in H, M, L p_H, p_M, p_L (sum to 1)

K-1 if K labels

Class conditional distribution of features:

$$\begin{matrix} \mu_1 \dots \mu_K & - & Kd \\ \Sigma_1 \dots \Sigma_K & - & Kd \left(\frac{d+1}{2} \right) \end{matrix}$$

$P(X=x | Y = y) \sim N(\mu_y, \Sigma_y)$ for each y

μ_y - d-dim vector

Σ_y - dxd matrix

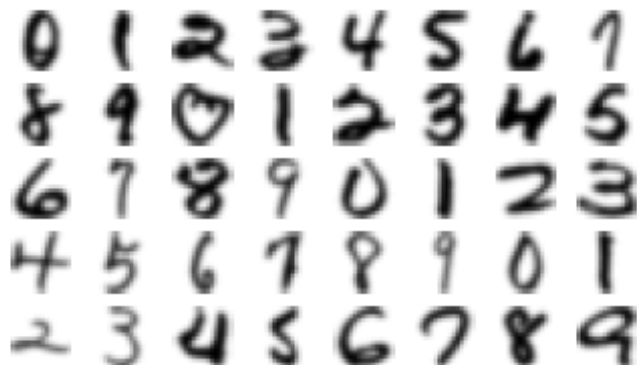
$Kd + Kd(d+1)/2 = O(Kd^2)$ if d features

$$3 \times (256^2)^2$$

Quadratic in dimension d! If d = 256x256 pixels, ~ 13 billion parameters!

Multi-class, multi-dimensional classification - Discrete features

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \{0,1\}^d$$



"0"
"1"
...
"9"

Input feature vector, X

Label, Y

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_{|V|} \end{bmatrix}$$



Sports
Science
News

Input feature vector, X

Label, Y

How many parameters do we need to learn (discrete features)?

Class probability: $K = 10$

$P(Y = y) = p_y$ for all y in $0, 1, 2, \dots, 9$ p_0, p_1, \dots, p_9 (sum to 1)

K-1 if K labels

Class conditional distribution of (binary) features:

$P(X=x | Y = y) \sim$ For each label y , maintain probability table with $2^d - 1$ entries

$\begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$

$X \in 2^d$

0, 0, ..., 0
0, 0, ..., 1
0, 0, ..., 10
⋮

$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$

$K(2^d - 1)$ if d binary features

Exponential in dimension d !

$x_i \in \{0, 1\}$

pixel

What's wrong with too many parameters?

- How many training data needed to learn one parameter (bias of a coin)?



- Need lots of training data to learn the parameters!
 - Training data $>$ number of (independent) parameters

runtime + storage requirement

Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:

- Features are independent given class:

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$

$$\begin{aligned} \underline{P(X_1, X_2|Y)} &= P(X_1|\underline{X_2}, \underline{Y})P(X_2|Y) \leftarrow \text{chain rule} \\ &= \underline{P(X_1|Y)}\underline{P(X_2|Y)} \end{aligned}$$

- More generally:

$$P(\underline{X_1 \dots X_d|Y}) = \prod_{i=1}^d P(X_i|Y)$$

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \dots \\ X_d \end{bmatrix}$$

- If conditional independence assumption holds, NB is optimal classifier! But worse otherwise.

Conditional Independence

- X is **conditionally independent** of Y given Z:

probability distribution governing X is independent of the value of Y, given the value of Z

$$(\forall x, y, z) P(\underbrace{X = x}_{\text{X}} | \underbrace{Y = y, Z = z}_{\text{Y, Z}}) = P(\underbrace{X = x}_{\text{X}} | \underbrace{Z = z}_{\text{Z}})$$

- Equivalent to:

$$P(\underbrace{X, Y}_{\text{X, Y}} | \underbrace{Z}_{\text{Z}}) = P(\underbrace{X}_{\text{X}} | \underbrace{Z}_{\text{Z}}) P(\underbrace{Y}_{\text{Y}} | \underbrace{Z}_{\text{Z}})$$

- e.g., $P(\text{Thunder} | \text{Rain}, \text{Lightning}) = P(\text{Thunder} | \text{Lightning})$

Note: does NOT mean Thunder is independent of Rain

Conditional vs. Marginal Independence

London taxi drivers: A survey has pointed out a positive and significant correlation between the number of accidents and wearing coats. They concluded that coats could hinder movements of drivers and be the cause of accidents. A new law was prepared to prohibit drivers from wearing coats when driving.

Finally another study pointed out that people wear coats when it rains...

Wearing coats is independent of accidents conditioning on the fact that it rained

Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:
 - Features are independent given class:

$$\underline{P(X_1 \dots X_d | Y)} = \prod_{i=1}^d P(X_i | Y)$$

$$\begin{aligned} f_{NB}(\mathbf{x}) &= \arg \max_y P(x_1, \dots, x_d | y) P(y) \\ &= \arg \max_y \prod_{i=1}^d \underline{P(x_i | y)} P(y) \end{aligned}$$

- How many parameters now?

How many parameters do we need to learn (continuous features)?

Class probability:

$$P(Y = y) = p_y \text{ for all } y \text{ in } H, M, L \quad p_H, p_M, p_L \text{ (sum to 1)}$$

K-1 if K labels

Class conditional distribution of features (using Naïve Bayes assumption):

$$P(\underline{X}_i = x_i | Y = y) \sim N(\underline{\mu}_i^{(y)}, \underline{\sigma}_i^{2(y)}) \text{ for each } y \text{ and each pixel } i$$

2Kd if d features

Linear instead of Quadratic in dimension d!

Handwritten diagram illustrating the parameter count for a class y . A bracket groups the parameters $\mu_1^{(y)}, \sigma_1^{2(y)}, \mu_2^{(y)}, \sigma_2^{2(y)}, \dots, \mu_d^{(y)}, \sigma_d^{2(y)}$. Below the bracket, the expression $Kd + \frac{Kd(d+1)}{2}$ is written, representing the total number of parameters for K classes and d features.

How many parameters do we need to learn (discrete features)?

Class probability:

$$P(Y = y) = p_y \text{ for all } y \text{ in } 0, 1, 2, \dots, 9 \quad p_0, p_1, \dots, p_9 \text{ (sum to 1)}$$

K-1 if K labels

Class conditional distribution of (binary) features:

$P(X_i = x_i | Y = y)$ – one probability value for each y , pixel i

Kd if d binary features

Linear instead of Exponential in dimension d!

$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix} \rightarrow \{0,1\}$$

$$P(X = x | Y = y) \\ = P(x_1 = x_1, \dots, x_d = x_d | Y = y)$$

Naïve Bayes Classifier

- Bayes Classifier with additional “naïve” assumption:
 - Features are independent given class:

$$P(X_1 \dots X_d | Y) = \prod_{i=1}^d P(X_i | Y)$$

$$\begin{aligned} f_{NB}(\mathbf{x}) &= \arg \max_y P(x_1, \dots, x_d | y) P(y) \\ &= \arg \max_y \prod_{i=1}^d P(x_i | y) P(y) \end{aligned}$$

- Has fewer parameters, and hence requires fewer training data, even though assumption may be violated in practice

Naïve Bayes Algo – Discrete features

- Training Data $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n$ $X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$

- Maximum Likelihood Estimates

- For Class probability

$$\hat{P}(y) = \frac{\{\#j : Y^{(j)} = y\}}{n} \leftarrow \text{data with label } y$$

- For class conditional distribution

$$\hat{P}(x_i | y) = \frac{\hat{P}(x_i, y)}{\hat{P}(y)} = \frac{\{\#j : X_i^{(j)} = x_i, Y^{(j)} = y\} / n}{\{\#j : Y^{(j)} = y\} / n} \leftarrow$$

- NB Prediction for test data $X = (x_1, \dots, x_d)$

$$Y = \arg \max_y \hat{P}(y) \prod_{i=1}^d \frac{\hat{P}(x_i, y)}{\hat{P}(y)}$$

Issues with Naïve Bayes

- **Issue 1:** Usually, features are not conditionally independent:

$$P(X_1 \dots X_d | Y) \neq \prod_i P(X_i | Y)$$

Nonetheless, NB is the single most used classifier particularly when data is limited, works well

- **Issue 2:** Typically use **MAP** estimates instead of MLE since insufficient data may cause MLE to be zero.

Insufficient data for MLE

- What if you never see a training instance where $X_1=a$ when $Y=b$?

– e.g., $b=\{\text{SpamEmail}\}$, $a = \{\text{‘Earn’}\}$

– $\hat{P}(X_1 = a \mid Y = b) = 0$

- Thus, no matter what the values X_2, \dots, X_d take:

$$\hat{P}(X_1 = a, X_2 \dots X_d \mid Y) = \hat{P}(X_1 = a \mid Y) \prod_{i=2}^d \hat{P}(X_i \mid Y) = 0$$

- What now???

Naïve Bayes Algo – Discrete features

- Training Data $\{(X^{(j)}, Y^{(j)})\}_{j=1}^n$ $X^{(j)} = (X_1^{(j)}, \dots, X_d^{(j)})$
- Maximum A Posteriori (MAP) Estimates – add m “virtual” data

Assume priors

$$Q(Y = \overset{\text{Spam}}{b})$$

$$Q(X_i = \overset{\text{Earn}}{a}, Y = \overset{\text{Spam}}{b})$$

$$\hat{P}(X_i = a | Y = b) = \frac{\{\#\text{j} : X_i^{(j)} = a, Y^{(j)} = b\} + mQ(X_i = a, Y = b)}{\{\#\text{j} : Y^{(j)} = b\} + mQ(Y = b)}$$

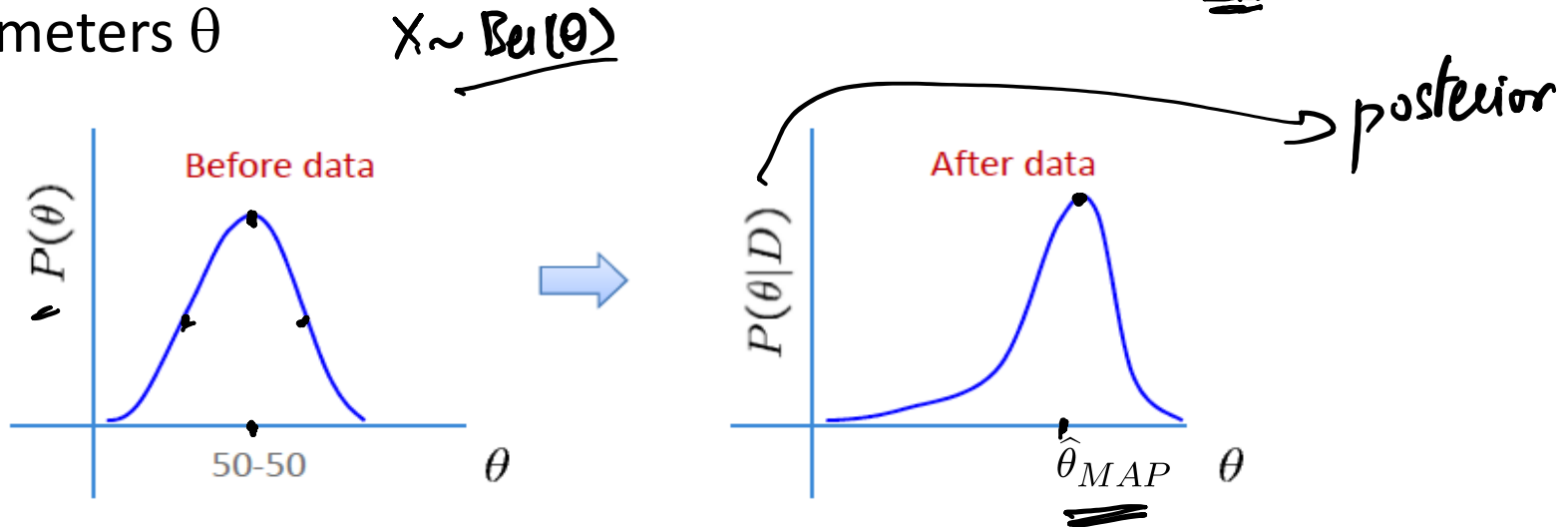
virtual examples
with $Y = b$

Now, even if you never observe a class/feature posterior probability never zero.

Max A Posteriori (MAP) estimation

Justification for adding virtual examples

- Assume a prior (before seeing data D) distribution $P(\theta)$ for parameters θ



- Choose value that maximizes a posterior distribution $P(\theta | D)$ of parameters θ

$$\begin{aligned} \hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} P(D | \theta) P(\theta) \end{aligned}$$

prior

How to choose prior distribution?

- $P(\theta)$



- Prior knowledge about domain e.g. unbiased coin $P(\theta) = \frac{1}{2}$

- A mathematically convenient form e.g. “conjugate” prior
 - If $P(\theta)$ is conjugate prior for $P(D|\theta)$, then Posterior has same form as prior

$$\text{Posterior} \equiv \text{Likelihood} \times \text{Prior}$$

$$\underline{P(\theta|D)} \equiv P(D|\theta) \times \underline{P(\theta)}$$

e.g.	[<u>Beta</u>	<u>Bernoulli</u>	<u>Beta</u>	$\theta = \text{bias}$
		<u>Gaussian</u>	<u>Gaussian</u>	<u>Gaussian</u>	$\theta = \text{mean } \underline{\mu}$ (known Σ)
	[<u>inv-Wishart</u>	<u>Gaussian</u>	<u>inv-Wishart</u>	$\theta = \text{cov matrix } \Sigma$ (known μ)

MAP estimation for Bernoulli r.v.

Choose θ that maximizes a posterior probability

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} \underbrace{P(D | \theta)} P(\theta)\end{aligned}$$

$X \sim \text{Beta}(\theta)$
 $\in \{0, 1\}$
 $\{H, T\}$

MAP estimate of probability of head (using Beta conjugate prior):

$$P(\theta) \sim \text{Beta}(\beta_H, \beta_T)$$

$$\theta = P(X=H)$$

Beta distribution

$Beta(\beta_H, \beta_T)$

More concentrated as values of β_H, β_T increase

