

Naïve Bayes contd...

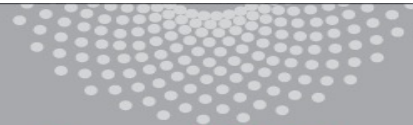
Aarti Singh

Machine Learning 10-315

Feb 2, 2022

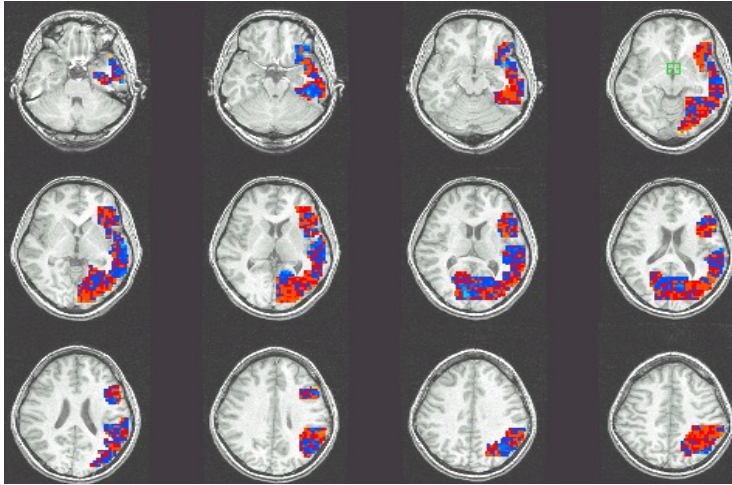


MACHINE LEARNING DEPARTMENT



Carnegie Mellon.
School of Computer Science

Multi-class, multi-dimensional classification – Continuous features



High Stress
Moderate Stress
Low Stress

Input feature vector, $X = \begin{bmatrix} x_{(1)} \\ \vdots \\ x_{(d)} \end{bmatrix}$ **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:

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

K-1 if K labels

Class conditional distribution of features:

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

μ_y - d-dim vector
 Σ_y - dxd matrix

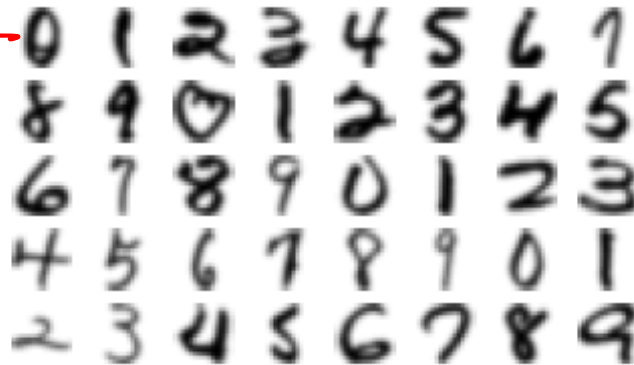
$$Kd + K \frac{d(d+1)}{2}$$

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

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

Multi-class, multi-dimensional classification - Discrete features

$$X = \begin{bmatrix} 0_{(1)} \\ \vdots \\ 0_{(d)} \end{bmatrix}$$



"0"
"1"
⋮
"9"

Input feature vector, X

Label, Y

$$X = \begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix}$$



Sports
Science
News

Input feature vector, X

Label, Y

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

Class probability:

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

K-1 if K labels

Class conditional distribution of (binary) features:

$P(X=x | Y = y) \sim$ For each label y , maintain probability table with

$\begin{bmatrix} 0/1 \\ \vdots \\ 0/1 \end{bmatrix}_d$ $\overset{=}{\sim}$ $\underline{2^d - 1}$ entries

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

Exponential in dimension d !

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

Naïve Bayes Classifier

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

$$X = \begin{bmatrix} X_{(1)} \\ X_{(2)} \end{bmatrix}$$

$$P(X|Y) = P(X_1, X_2|Y) = P(X_1|X_2, Y)P(X_2|Y) \\ = P(X_1|Y) \cdot P(X_2|Y)$$

$P(A,B) = P(A|B)P(B)$

- More generally:

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

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

$$X = \begin{bmatrix} X_1 \\ \vdots \\ X_d \end{bmatrix}$$

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

Conditional Independence

$$P(X, Y | Z) = P(X | Z)P(Y | Z) \checkmark$$

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

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

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

- Equivalent to:

$$P(X, Y | Z) = P(X | Z)P(Y | Z) \checkmark$$

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

Note: does NOT mean Thunder is independent of Rain

$Y = \text{low stress}$) $X_1 \dots X_n$ $\begin{bmatrix} X_{(1)} \\ \vdots \\ X_{(d)} \end{bmatrix}$

Naïve Bayes Classifier

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

$\sim N(\mu_y, \Sigma_y)$

$$P(\mathbf{X}|Y) = P(X_{(1)} \dots X_{(d)}|Y) = \prod_{i=1}^d P(X_{(i)}|Y)$$

$d \times d$

$$\Sigma_y = \begin{bmatrix} \sigma_{(1)}^2 & 0 \\ 0 & \sigma_{(d)}^2 \end{bmatrix}$$

$\Sigma_y(i,j) = E[(X_i - E(X_i))(X_j - E(X_j))]$

Y

$$f_{NB}(\mathbf{x}) = \arg \max_y P(x_{(1)}, \dots, x_{(d)} | y) P(y)$$

$$= \arg \max_y \prod_{i=1}^d P(x_{(i)} | y) P(y)$$

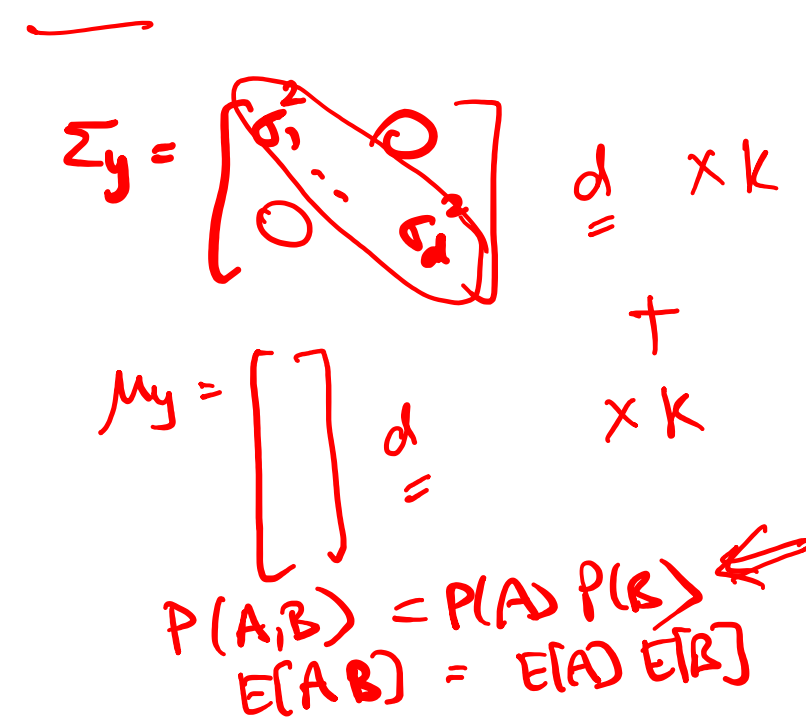
- How many parameters now?

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

$$E[(x_i - E x_i) | Y] = 0$$

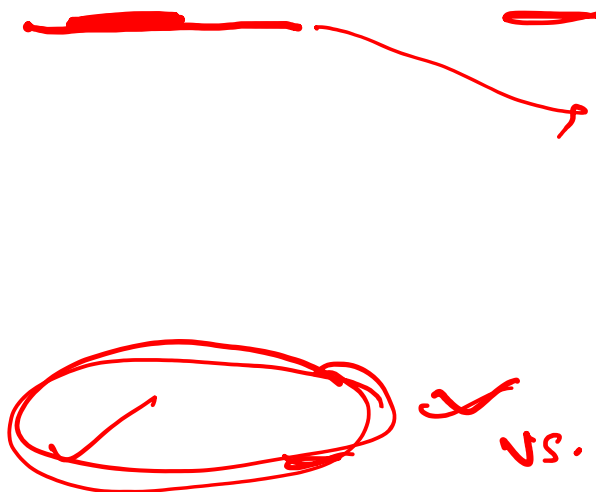
➤ Poll

$$P(X|Y) \sim \mathcal{N}(\mu_y, \Sigma_y) \quad (\Sigma_y)_{ij} = \underline{E[(X_i - E X_i)(X_j - E X_j) | Y]}$$



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

➤ Poll


$$\begin{aligned} & P(X_1=x_1, X_2=x_2, \dots, X_d=x_d | Y) \\ &= \prod_{i=1}^d \underbrace{P(X_i=x_i | Y)}_d \end{aligned}$$

vs. $O(K^{2^d})$

Naïve Bayes Classifier

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

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

$$f_{NB}(\mathbf{x}) = \arg \max_y P(x_1, \dots, x_d | y) P(y)$$

$$= \arg \max_y \prod_{i=1}^d P(x_i | y) P(y)$$

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

Learned Gaussian Naïve Bayes Model Means for $P(\text{BrainActivity} \mid \text{WordCategory})$

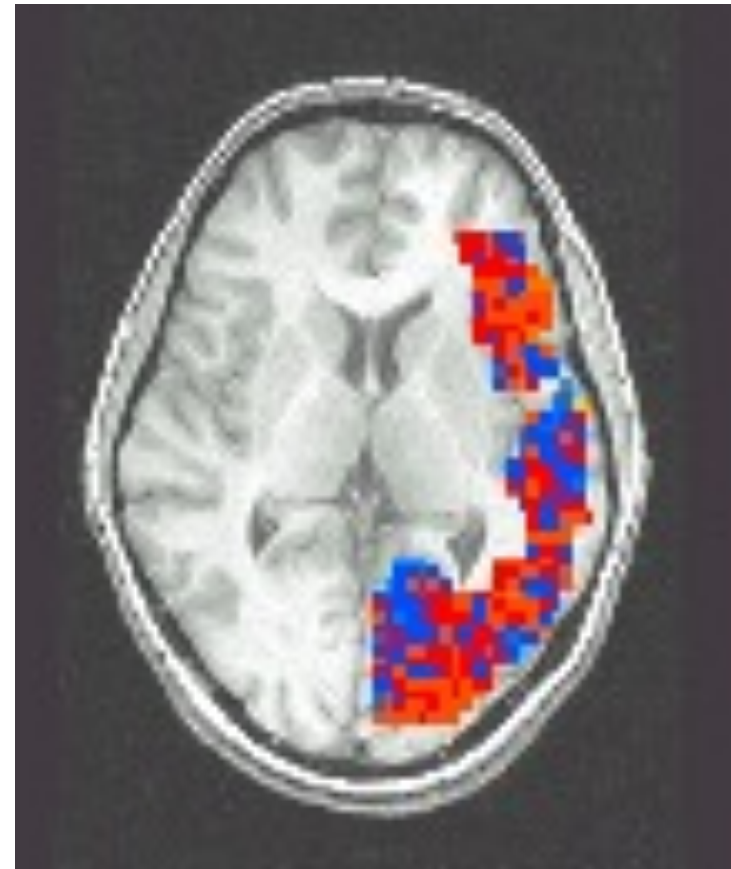
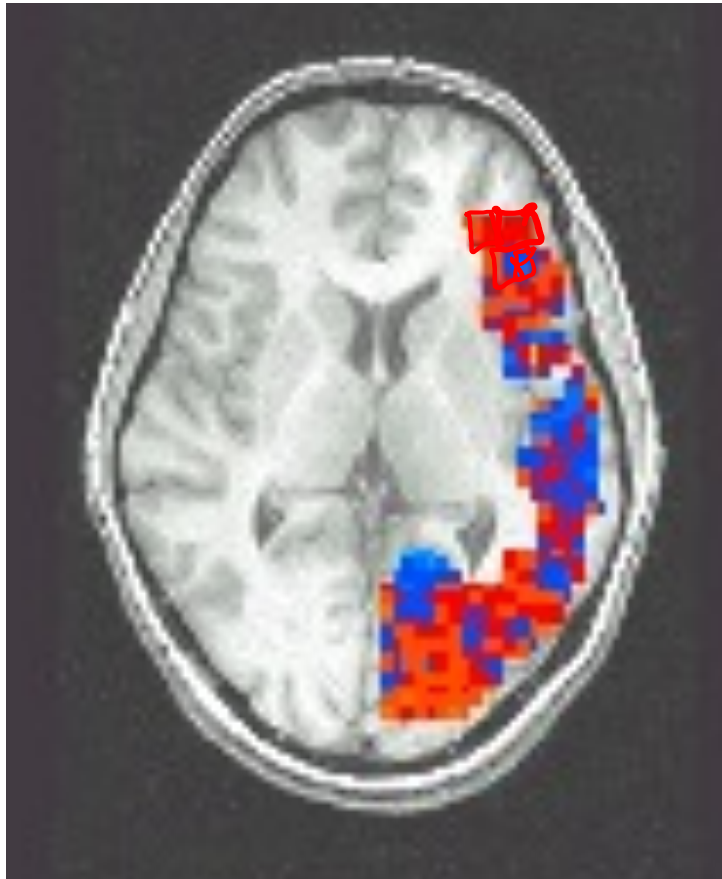
Pairwise classification accuracy: **85%** [Mitchell et al.03]

μ_y Σ_y

People words



Animal words



Text classification

Raw input, ~~X~~



Features

$$\begin{bmatrix} X_{(1)} \\ \vdots \\ X_{(d)} \end{bmatrix}$$



Model for input features



| | |
|-------|----|
| word1 | 5 |
| word2 | 2 |
| word3 | 10 |
| word4 | 20 |
| word5 | 12 |
| word6 | 5 |
| word7 | 8 |
| word8 | 4 |
| . | . |
| . | . |
| . | . |

Bag of words

$$P(X=x | Y=y) = P(\text{word1} = 5, \text{word2} = 2, \text{word3} = 10, \dots | Y=y)$$

$$= \prod_{i=1}^{|\text{word}|} P(X_{(i)} = x_{(i)} | Y=y)$$

HW1!

Bag of words + Naïve Bayes