

Consider the following training data points and their corresponding  $\alpha$  values

$$\begin{aligned} \mathbf{x}^1 &= \begin{bmatrix} 2 \\ 0 \\ -1.5 \end{bmatrix}, y^1 = +1, \alpha^1 = 0.5 & \mathbf{x}^2 &= \begin{bmatrix} 1 \\ -0.5 \\ 0 \end{bmatrix}, y^2 = +1, \alpha^2 = 2 \\ \mathbf{x}^3 &= \begin{bmatrix} 5 \\ 0 \\ 2 \end{bmatrix}, y^3 = -1, \alpha^3 = 2 & \mathbf{x}^4 &= \begin{bmatrix} 13 \\ -8 \\ 7 \end{bmatrix}, y^4 = -1, \alpha^4 = 0 \\ \mathbf{x}^5 &= \begin{bmatrix} 7 \\ -6 \\ 0 \end{bmatrix}, y^5 = -1, \alpha^5 = 0.5 & \mathbf{x}^6 &= \begin{bmatrix} -5 \\ 0 \\ -5 \end{bmatrix}, y^6 = +1, \alpha^6 = 0 \end{aligned}$$

Which of the above are considered support vectors?

Use the above data to compute the separating hyperplane  $\mathbf{w}$  found by the linear SVM.

Consider a Bayesian classifier and a data set of containing 30 features. Assume there are 100 training data points. We report the highest Likelihood classifier given  $n$  features below:

#Feats	1	2	3	4	5	6	7	...	21
<b>MaxL</b>	$6 \times 10^{-23}$	$1 \times 10^{-21}$	$2 \times 10^{-20}$	$4 \times 10^{-20}$	$8 \times 10^{-19}$	$2 \times 10^{-19}$	$4 \times 10^{-19}$		$1 \times 10^{-9}$

#Feats	22	23	24	25	26	27	28	29	30
<b>MaxL</b>	$1 \times 10^{-8}$	$1 \times 10^{-7}$	$6 \times 10^{-7}$	$4 \times 10^{-6}$	$1 \times 10^{-5}$	$5 \times 10^{-5}$	$9 \times 10^{-5}$	$2 \times 10^{-4}$	$3 \times 10^{-4}$

What is the first feature removal step (how many features left?) when AIC no longer improves?

What is the first feature selection step (how many features chosen?) when BIC no longer improves?

Consider the following Principal Components:

$$\mathbf{u}_1 = \begin{bmatrix} 0.41 \\ -0.41 \\ 0 \\ 0.82 \end{bmatrix} \quad \mathbf{u}_2 = \begin{bmatrix} 0 \\ 0.53 \\ 0.80 \\ 0.27 \end{bmatrix}$$

For each data point below, calculate the weights  $z$  for each principal component.

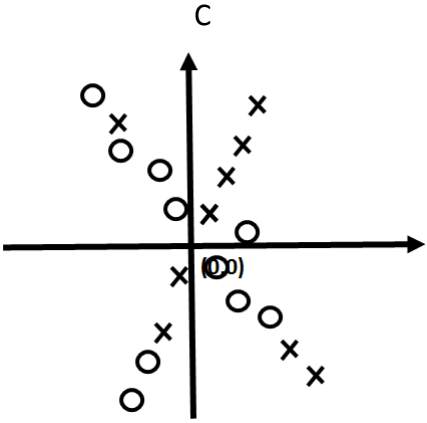
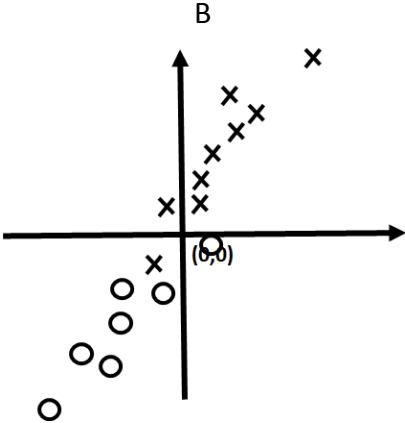
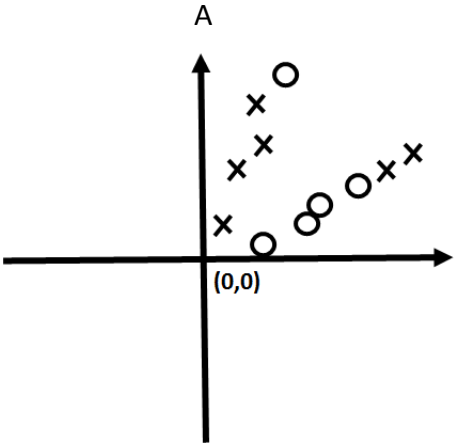
$$\mathbf{x}^1 = \begin{bmatrix} -4.4 \\ 3.8 \\ -1.0 \\ -9.2 \end{bmatrix}$$

$$\mathbf{x}^2 = \begin{bmatrix} 0.1 \\ 0.7 \\ 1.2 \\ 0.7 \end{bmatrix}$$

$$\mathbf{x}^3 = \begin{bmatrix} 2.3 \\ -2.5 \\ -0.3 \\ 4.4 \end{bmatrix}$$

Using only  $\mathbf{u}_1$ , compute the reconstruction error for each of the above data points.

We observe the following data points in two dimensions. We wish to define new dimensions to better describe the data. Should we use PCA, ICA, or NMF?



Consider an HMM with 3 states:

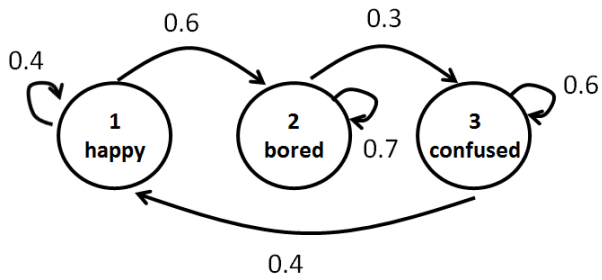
1: happy, 2: bored, 3: confused

And 4 potential outputs at each time (each output is punctuation mark)

1: ! 2: ? 3: . 4: -

The transition matrix is capture by the diagram

The initial probabilities are:



$$\begin{aligned} \pi_{happy} &= 0.2 \\ \pi_{bored} &= 0.6 \\ \pi_{confused} &= 0.2 \end{aligned}$$

The emission matrix is  $\phi$ :

State\Obs	!	?	.	-
Happy	0.8	0.1	0.1	0
Bored	0	0.2	0.7	0.1
Confused	0.1	0.7	0	0.2

Which of the following are impossible state sequences:

- happy, confused, happy, happy
- bored, bored, bored, confused
- happy, happy, bored, confused, happy
- bored, bored, happy, happy

Presume the sequence:

?, ?, ?, ?

What is the forward value  $\alpha_1(\text{happy})$

What is the forward value  $\alpha_2(\text{confused})$

**Note: Viterbi was only briefly covered in class. I include this problem for review/clarification of the process but will make Viterbi worth few points (if any) on the exam.**

Using the Viterbi algorithm, what is the sequence of most probable states for each observation sequence below:

!, ?, !

How can we compute the emission probabilities for each state if we have a set of training sequences where both the observations and the underlying states are known? (Describe in a sentence and/or write a mathematical expression.)

How can we compute the emission probabilities if we have a set of training sequences where the observations are known but the underlying states are NOT known? (Describe in 1-3 sentences and/or write 1-3 mathematical expressions.)

We seek to apply MAP to learn  $P(X|Y; \theta)$ , where  $X$  is a binary feature “likes peanuts?” (yes or no) and  $Y$  represents the class “is an elephant” (true or false). We use  $\alpha_{true}$  and  $\beta_{true}$  to represent the probability prior belief  $X=yes$  and  $X=no$  given  $Y=true$ . **Match each  $\alpha_{true}$ ,  $\beta_{true}$  pairing on the left with its corresponding meaning on the right.**

$$\alpha_{true} = 1000, \beta_{true} = 1000$$

Strong belief elephants don't like peanuts

$$\alpha_{true} = 1, \beta_{true} = 0.2$$

Mild belief elephants like peanuts

$$\alpha_{true} = 1, \beta_{true} = 2000$$

Strong belief elephants are equally likely to like or dislike peanuts

What are the effects of changing the following parameters in gradient ascent learning for logistic regression?

1:  $\epsilon$

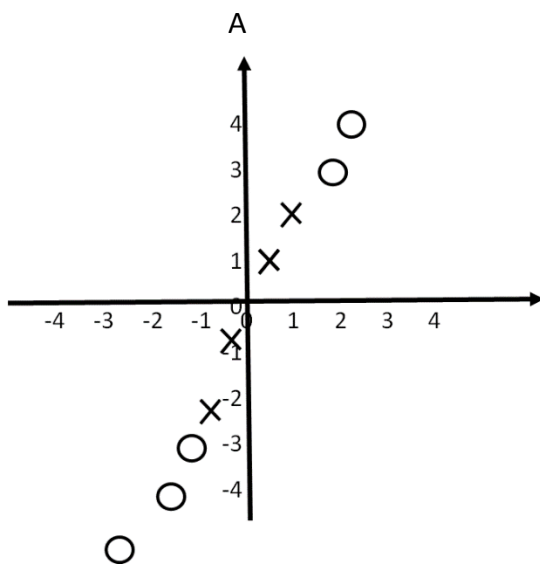
2:  $\lambda$

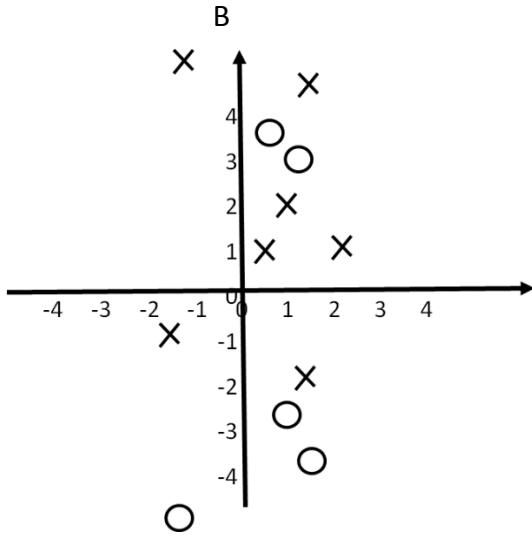
We have data points with 10 features each. Each feature is a numeric value along the real number line. We wish to learn a classifier to label each data point as one of four classes. Using logistic regression (including +b term!), how many parameters must we learn?

How many parameters must we learn if we wish to classify the data above into only **two** classes using a linear SVM?

How many parameters must we learn if we wish to classify the data above into only two classes using a quadratic kernel SVM?

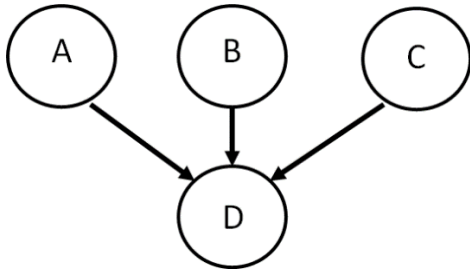
Define a mapping function that will allow a linear separator to distinguish between the two classes for each set of data points. ( $x_1$  is horizontal axis,  $x_2$  is vertical axis)



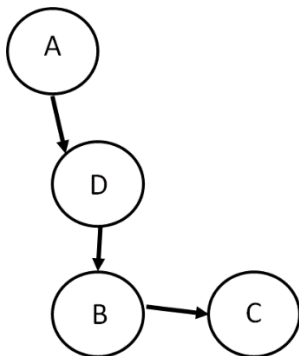


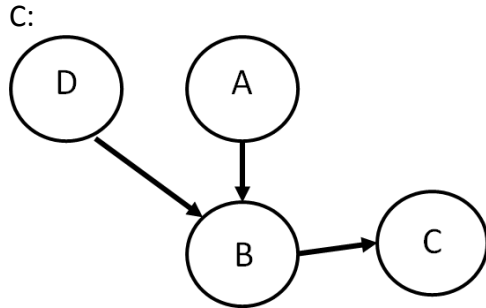
Write the formula to find the total joint probability for each of the following Bayesian Networks (to find  $P(A,B,C,D)$ ).

A:



B:





Compute  $P(\text{subset of variables})$

For each Bayes net above, how can we find the value of  $P(B)$

Let us presume we wish to classify whether an animal is a mammal or a bird, and we use 30 features  $x_1, \dots, x_{30}$  as a basis for classification. (For example,  $x_1$  can be size,  $x_2$  can be typical speed of motion,  $x_3$  can be blood temperature, etc.) For a given animal we observe the

following feature vector:  $x = \begin{bmatrix} 10 \\ 1 \\ 55 \\ \vdots \\ 2 \\ 16 \\ 240 \end{bmatrix}$

and we are informed of the following probabilities:

$P(x_1=10 | Y=\text{bird})=0.01$        $P(x_2=1 | Y=\text{bird})=0.2$     ...     $P(x_{29}=16 | Y=\text{bird})=0.05$   
 $P(x_{30}=240 | Y=\text{bird})=0.2$

$P(x_1=10 | Y=\text{mammal})=0.06$      $P(x_2=1 | Y=\text{mammal})=0.08$  ...  $P(x_{29}=16 | Y=\text{mammal})=0.1$   
 $P(x_{30}=240 | Y=\text{mammal})=0.07$

If we use Naïve Bayes on a standard computer, we will find  $P(x_1, \dots, x_{30} | Y=\text{bird})=0$  and  $P(x_1, \dots, x_{30} | Y=\text{mammal})=0$ , despite the fact that **none** of the terms  $P(x_i | Y=\text{mammal})=0$  nor  $P(x_i | Y=\text{bird})=0$

a) Why does this happen on a computer?

b) What mathematical operation can we use to prevent this problem?