

INF8953 CE - Fall 2020

# Machine Learning

- Sarath Chandar

J. Probabilistic Generative Models .



## F. Probabilistic Generative Models

### Generalized linear models:-

In linear regression, the model prediction  $y(n; w)$  was given by a linear function of the parameters ' $w$ '.

In the simplest case, the model is also linear in the input variables.

$$y(n) = \bar{w}^T x + w_0$$

$\uparrow$   
real number.

For classification, we wish to predict discrete class labels or posterior probabilities that lie in the range  $(0, 1)$ .

↳ We transform the linear function of  $w$  using a non-linear function  $f(\cdot)$  so that

$$y(n) = f(w^T x + w_0)$$

$f(\cdot)$  - activation function.

Note: This model is no longer linear in the parameters due to the presence of non-linear  $f(\cdot)$ .

The decision surface is given by

$$y(n) = \text{constant}$$

proving that the decision surface is linear.

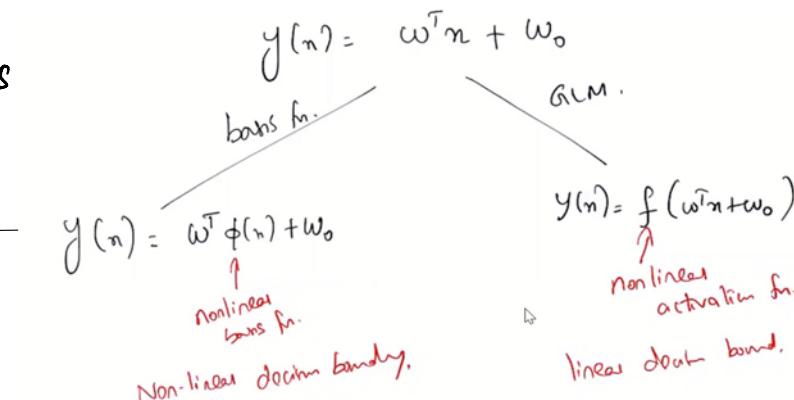
$$f(w^T n + w_0) = \text{constant}$$

$$w^T n + w_0 = f^{-1}(\text{constant})$$

$$w^T n + w_0 = \text{constant}$$

Decision surfaces are linear functions of  $x$ , even if  $f(\cdot)$  is non-linear! This type of models are called generalized linear models (GLM).

GLMs have complex analytical and computational properties than linear models. However, they are still relatively simple when compared to more general non-linear models



## Probabilistic Generative Models:-

Model the class-conditional densities

$P(x|C_k)$  and class priors  $P(C_k)$  and then use Bayes' theorem to complete posterior probabilities  $P(C_k|x)$ .

### 2-class problem:-

$$P(C_1|x) = \frac{P(x|C_1) P(C_1)}{P(x)}$$

$$= \frac{P(x|C_1) P(C_1)}{P(x|C_1) P(C_1) + P(x|C_2) P(C_2)}$$

mul and div by numerator

$$= \frac{1}{1 + \frac{P(x|C_2) P(C_2)}{P(x|C_1) P(C_1)}}$$

$$= \frac{1}{1 + \exp\left(-\ln \frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)}\right)}$$

$$= \frac{1}{1 + \exp(-a)} = \sigma(a)$$

where  $a = \ln \left( \frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)} \right)$

$$\sigma(a) = \frac{1}{1+e^{-a}}$$

(logistic Sigmoid fn).

Sigmoid function:-

→ Squashing function

→ range (0,1)

→ Symmetry property :  $\sigma(-a) = 1 - \sigma(a)$

→ inverse :  $a = \ln\left(\frac{\sigma}{1-\sigma}\right)$

↗  
logit function

logit function - log of the ratio of probabilities

$\ln\left(\frac{P(C_1|x)}{P(C_2|x)}\right)$  for the two classes, also known  
as the log odds.

For  $k > 2$  classes:

$$P(C_k|x) = \frac{P(x|C_k) P(C_k)}{\sum_j P(x|C_j) P(C_j)}$$

$$P(C_k | x) = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where  $a_k = \ln(P(x|C_k) P(C_k))$

Normalized exponential

Multiclass generalization of the logistic sigmoid.

also known as softmax function, as

it represents a smoothed version of the 'max' function because if  $a_k \gg a_j$  for all  $j \neq k$ , then  $P(C_k | x) \approx 1$  and  $P(C_j | x) \approx 0$ .

## Gaussian Discriminant Analysis:-

### Assumptions:-

① Class conditional densities are Gaussian.

② All classes share the same covariance matrix.

$$P(x|C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) \right\}$$

Mean :  $\mu_k$

Covariance :  $\Sigma$

Note:  $|\Sigma|$  is the determinant of  $\Sigma$ .

Consider 2-class problem:

$$P(C_1|x) = \sigma(a)$$

$$\text{where } a = \ln \left( \frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)} \right)$$

$$a = \ln \left( \frac{P(x|C_1)}{P(x|C_2)} \right) + \ln \left( \frac{P(C_1)}{P(C_2)} \right)$$

$$= \cancel{\ln} \exp \left\{ -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)^T \Sigma^{-1} (x - \mu_2) \right\}$$

$$+ \ln \frac{P(C_1)}{P(C_2)}$$

$$= -\frac{1}{2} \left( \cancel{x^T \Sigma^{-1} x} - x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 \right. \\ \left. - \cancel{x^T \Sigma^{-1} x} + x^T \Sigma^{-1} \mu_2 + \mu_2^T \Sigma^{-1} x - \mu_2^T \Sigma^{-1} \mu_2 \right) \\ + \ln \left( \frac{P(C_1)}{P(C_2)} \right)$$

$$\begin{aligned}
&= \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_1 - \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_2 - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 \\
&\quad + \ln \frac{P(C_1)}{P(C_2)} \\
&= \mathbf{x}^T (\bar{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)) - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 \\
&\quad + \ln \frac{P(C_1)}{P(C_2)} \\
&= \mathbf{w}^T \mathbf{x} + w_0
\end{aligned}$$

where  $\mathbf{w} = \bar{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{P(C_1)}{P(C_2)}$$

$$P(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

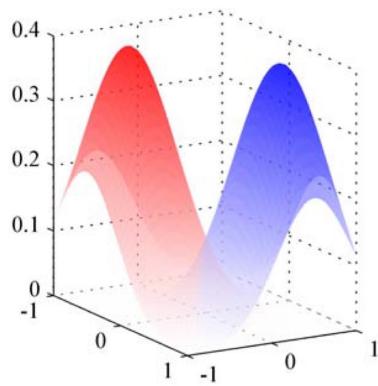
Note: Quadratic terms in  $\mathbf{x}$  from the exponents of the Gaussian densities have cancelled due to the assumption of the common covariance matrices.

The model is a linear function of  $\mathbf{x}$ .

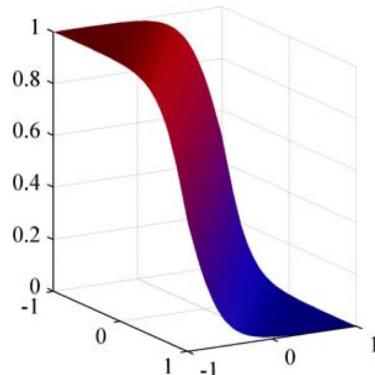


decision boundary is also a linear function of  $\mathbf{x}$ .

Prior probabilities  $P(C_k)$  enter only through the bias parameter  $w_0$  so that the changes in the prior have the effect of making parallel shifts of the decision boundary.



Class Conditional density  
for 2 classes.



$$P(c_1|x)$$

Parameters of the model:  $\mu_1, \mu_2, \Sigma, P(c_1), P(c_2)$

How to learn the parameters?

Consider a dataset  $\mathcal{D} = \{x^{(n)}, t^{(n)}\}_{n=1}^N$  of  $N$  examples.

$$t_n = 1 \text{ class } c_1$$

$$t_n = 0 \text{ class } c_2$$

Let  $P(C_1) = \pi$  Prior probability for  $C_1$

$$P(C_2) = 1 - \pi$$

Consider data point  $x$ ,

$$P(x^{(n)}, C_1) = P(C_1) P(x^{(n)} | C_1) = \pi \cdot N(x^{(n)} | \mu_1, \Sigma)$$

$$P(x^{(n)}, C_2) = P(C_2) P(x^{(n)} | C_2) = (1 - \pi) N(x^{(n)} | \mu_2, \Sigma)$$

Likelihood of the datapoint  $x^{(n)}, t^{(n)}$  is given by

$$\left[ \pi N(x^{(n)} | \mu_1, \Sigma) \right]^{t^{(n)}} \left[ (1 - \pi) N(x^{(n)} | \mu_2, \Sigma) \right]^{1 - t^{(n)}}$$

Assuming that the examples are independent and identically distributed (i.i.d)

$$P(D | \pi, \mu_1, \mu_2, \Sigma) = \prod_{n=1}^N \left[ \pi N(x^{(n)} | \mu_1, \Sigma) \right]^{t^{(n)}} \left[ (1 - \pi) N(x^{(n)} | \mu_2, \Sigma) \right]^{1 - t^{(n)}}$$

likelihood function.

We want to find the parameters that maximizes the likelihood function.

Maximum likelihood solution :-

We will often maximize log of the likelihood function instead of the likelihood fn. itself.

$$\begin{aligned} \ln P(D | \pi, \mu, \mu_2, \Sigma) &= \sum_{n=1}^N \ln \left[ \pi N(x^{(n)} | \mu, \Sigma) \right]^{\epsilon^{(n)}} \\ &\quad \left[ (1-\pi) N(x^{(n)} | \mu_2, \Sigma) \right]^{1-\epsilon^{(n)}} \\ &= \sum_{n=1}^N \left[ \epsilon^{(n)} \ln \pi N(x^{(n)} | \mu, \Sigma) + (1-\epsilon^{(n)}) \ln (1-\pi) N(x^{(n)} | \mu_2, \Sigma) \right] \end{aligned}$$

Maximization w.r.t.  $\pi$  :

$$\sum_{n=1}^N \left[ \epsilon^{(n)} \ln \pi + (1-\epsilon^{(n)}) \ln (1-\pi) \right] + \text{const}$$

diff. w.r.t.  $\pi$  and set it to zero.

$$\sum_{n=1}^N \left( \frac{\epsilon^{(n)}}{\pi} + \frac{1-\epsilon^{(n)}}{1-\pi} (-1) \right) = 0$$

$$\sum_{n=1}^N \left( \frac{t^{(n)}(1-\pi) - (1-t^{(n)})\pi}{\pi(1-\pi)} \right) = 0$$

$$\sum_{n=1}^N \left( t^{(n)} - \cancel{\frac{t^{(n)}}{\pi}} - \cancel{\pi} + \cancel{\frac{t^{(n)}}{\pi}} \right) = 0$$

$$\left( \sum_{n=1}^N t^{(n)} \right) - N\pi = 0$$

fraction of  
points  
in  $C_1$ .

$$\pi = \frac{\sum_{n=1}^N t^{(n)}}{N} = \frac{N_1}{N} = \frac{N_1}{N_1+N_2}$$

where  $N_1$  - number of data points in class  $C_1$ ,

$N_2$  - number of data points in class  $C_2$ .

Maximization w.r.t.  $\mu_1$ :

$$\sum_{n=1}^N t^{(n)} \ln N(x^{(n)} | \mu_1, \Sigma) + \text{Const.}$$

$$= -\frac{1}{2} \sum_{n=1}^N t^{(n)} (x^{(n)} - \mu_1)^T \Sigma^{-1} (x^{(n)} - \mu_1) + \text{Const}$$

diff. w.r.t.  $\mu_1$  and set it to zero.

$$\sum_{n=1}^N t^{(n)} (x^{(n)} - \mu_1) = 0$$

$$N_1 \mu_1 = \sum_{n=1}^N t^{(n)} x^{(n)}$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t^{(n)} x^{(n)}$$



mean of all the input vectors

$x^{(n)}$  assigned to class  $C_1$ .

$$\text{Similarly, } \mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1 - t^{(n)}) x^{(n)}$$

Maximization w.r.t.  $\Sigma$ :

Maximization w.r.t  $\Sigma$  is bit involved.

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$\text{where } S_1 = \frac{1}{N_1} \sum_{n \in C_1} (x^{(n)} - \mu_1)(x^{(n)} - \mu_1)^T$$

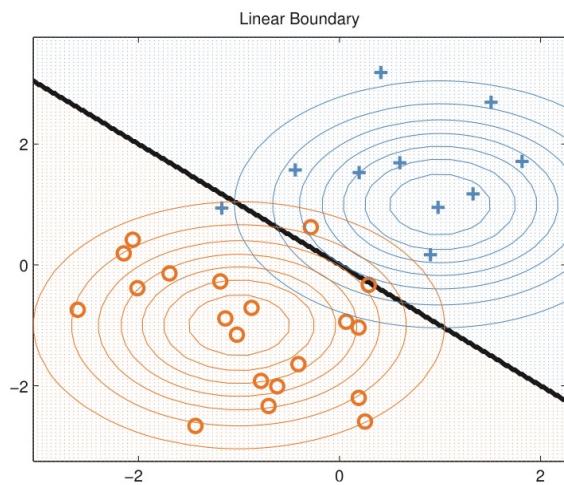
$$S_2 = \frac{1}{N_2} \sum_{n \in C_2} (x^{(n)} - \mu_2)(x^{(n)} - \mu_2)^T$$

Note: If we relax the assumption of a shared covariance matrix and allow each class conditional density  $p(x|C_k)$  to have its own

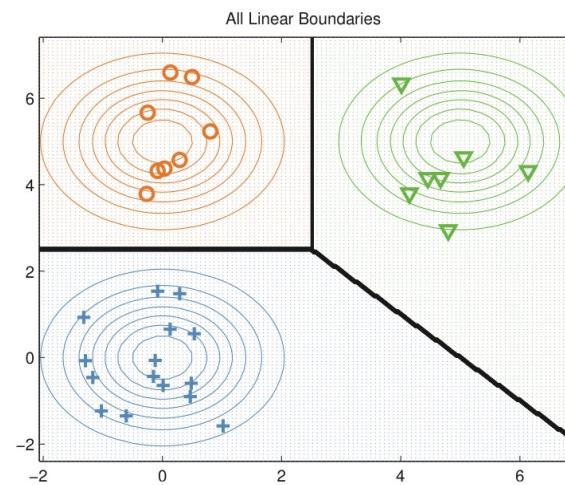
Covariance matrix  $\Sigma_k$ , then we will obtain quadratic functions of  $X$ .

→ Quadratic Discriminant Analysis (QDA)

LDA

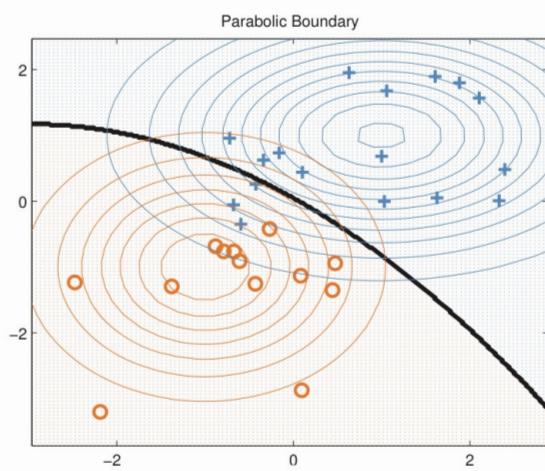


2 - classes

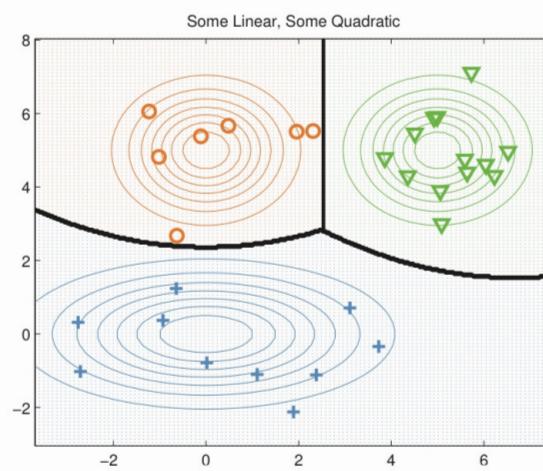


3 - classes

QDA



2 - classes



3 - classes

## Summary:

### Generative models for classification:-

Compute  $P(x|C_k)$  and  $P(C_k)$  and then use Bayes' law to compute  $P(C_k|x)$ .

GDA: Class conditional densities are Gaussians with shared Covariance matrix.

QDA: class conditional densities are Gaussians with separate covariance matrix.

If there are M features, then

GDA:  $kM$  parameters for the means,

$\frac{M(M+1)}{2}$  Parameters for the shared Covariance matrix.

$k-1$  parameters for  $P(C_k)$

QDA:  $kM$  parameters for means.

$\frac{k \cdot M \cdot (M+1)}{2}$  Parameters for Covariance matrices.

$k-1$  parameters for  $P(C_k)$ .

SQDA has much more parameters. But the shared Covariance assumption is too restrictive.

---

### Naive Bayes assumption:-

Features are conditionally independent given the class label.

$$\begin{aligned} P(x|C_k) &= P(x_1 \dots x_M | C_k) \\ &= \prod_{i=1}^M P(x_i | C_k) \end{aligned}$$

Now the Covariance matrix is diagonal.

This model is known as Gaussian Naive Bayes. (GNB)

GNB requires  $kM$  parameters for means,

$kM$  parameters for variance, and

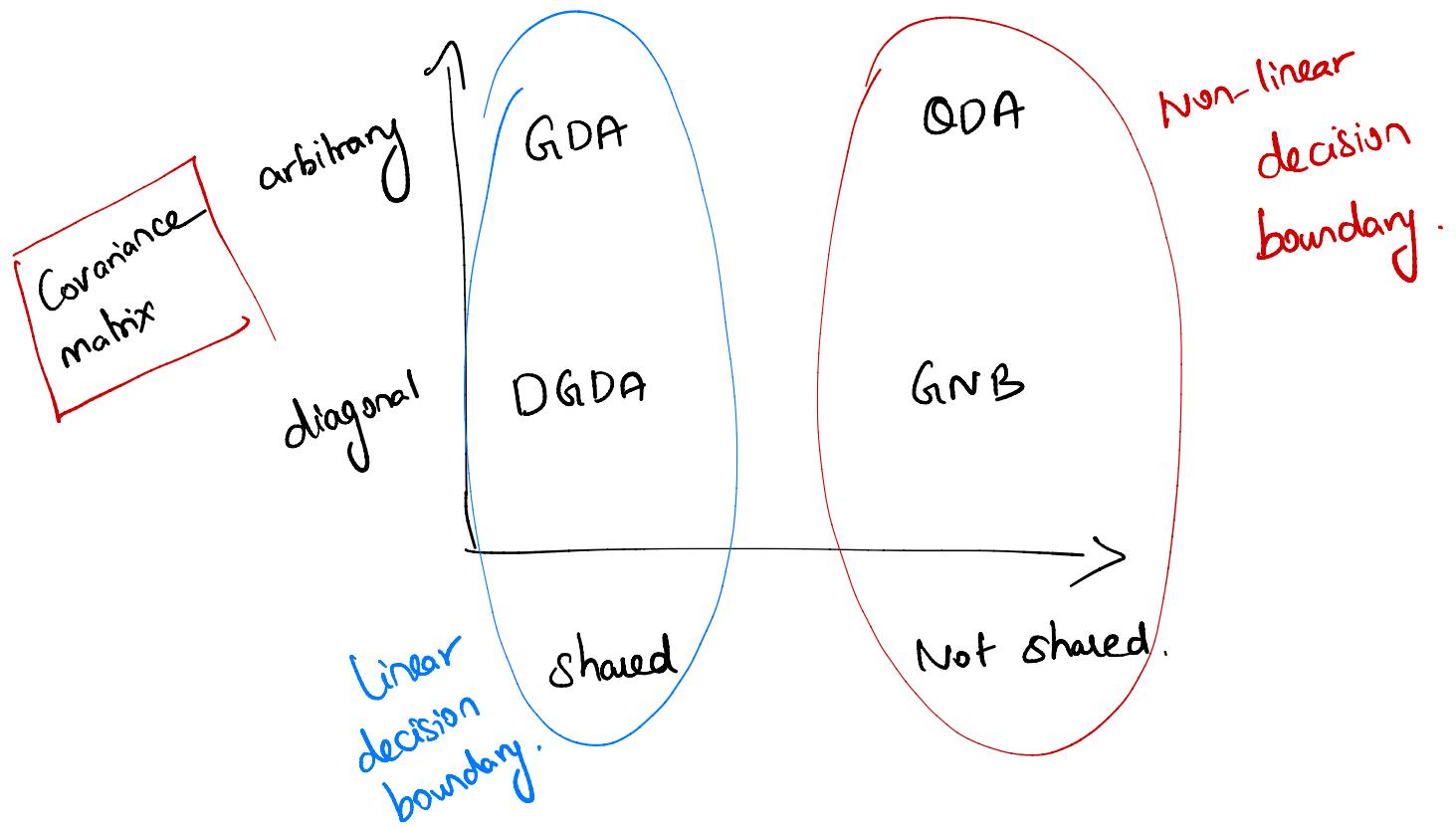
$K-1$  Parameters for  $P(C_k)$ .

This requires lesser number of parameters and still has non-linear decision boundary like QDA.

Note 1: If the diagonal Covariance matrix is shared among the classes, then the decision boundary will be linear. GNB with shared diagonal Covariance matrices is known as DGDA (Diagonal GDA).

Note 2: We can learn the Covariance matrix in all these methods by using maximum likelihood.

Summary: Continuous i/p features



Discrete features:-

Consider binary features:  $x_i \in \{0, 1\}$

If there are  $M$  inputs, then a general distribution would correspond to a table of  $2^M$  numbers for each class, containing  $2^{M-1}$  independent variables.

Naive Bayes' assumption :-  $P(x|C_k) = \prod_{i=1}^M P(x_i|C_k)$

$$P(x|C_k) = P(x_1 \dots x_M | C_k)$$

$$= P(x_1|C_k) P(x_2|C_k, x_1) P(x_3|C_k, x_1, x_2)$$

$$\dots P(x_M|C_k, x_1 \dots x_{M-1})$$

(chain rule)

$$= P(x_1|C_k) P(x_2|C_k) \dots P(x_M|C_k)$$

(Naive Bayes  
assumption)

$$= \prod_{i=1}^M P(x_i|C_k)$$

Each  $P(x_i|C_k)$  can be modeled as a Bernoulli with parameter  $\mu_{k_i}$  (i.e.  $P(x_i=1|C_k)$ )

$$P(x|C_k) = \prod_{i=1}^M \mu_{k_i}^{x_i} (1-\mu_{k_i})^{1-x_i}$$

This requires only  $M$  parameters for each class.

Consider 2 class problem :-

Likelihood :  $P(D|\mu_1, \mu_2, \pi)$

$$= \frac{N}{\prod_{i=1}^n} \left( \pi \prod_{i=1}^M \mu_{1i}^{x_i^{(n)}} (1-\mu_{1i}^{(n)}) \right)^{\pi} \left( (1-\pi) \prod_{i=1}^M \mu_{2i}^{x_i^{(n)}} (1-\mu_{2i}^{(n)}) \right)^{1-\pi}$$

$x$                                      $y$

$$\ln P(D | \mu_1, \mu_2, \pi) = \sum_{i=1}^N \left[ t^{(n)} \ln x_i + (1-t^{(n)}) \ln y_i \right]$$

$$\frac{\partial \ln P(D | \cdot)}{\partial \mu_{1,i}} = \frac{\partial}{\partial \mu_{1,i}} \left[ \sum_{n=1}^N \left( t^{(n)} \ln \mu_{1,i}^{x_i^{(n)}} (1-\mu_{1,i})^{1-x_i^{(n)}} \right) \right]$$

$$= \frac{\partial}{\partial \mu_{1,i}} \left[ \sum_{n=1}^N \left( t^{(n)} x_i^{(n)} \ln \mu_{1,i} + t^{(n)} (1-x_i^{(n)}) \ln (1-\mu_{1,i}) \right) \right]$$

diff. and set it to zero.

$$\sum_{n=1}^N \left( t^{(n)} x_i^{(n)} \frac{1}{\mu_{1,i}} + t^{(n)} (1-x_i^{(n)}) \frac{1}{(1-\mu_{1,i})} (-1) \right) = 0$$

$$\sum_{n=1}^N \left( t^{(n)} x_i^{(n)} (1-\mu_{1,i}) - t^{(n)} (1-x_i^{(n)}) \mu_{1,i} \right) = 0$$

$$\sum_{n=1}^N \left( t^{(n)} x_i^{(n)} - \cancel{t^{(n)} x_i^{(n)} \mu_{1,i}} - t^{(n)} \mu_{1,i} + \cancel{t^{(n)} x_i^{(n)} \mu_{1,i}} \right) = 0$$

$$\sum_{n=1}^N t^{(n)} \mu_{1,i} = \sum_{n=1}^N t^{(n)} x_i^{(n)}$$

$$\boxed{\mu_{1,i} = \frac{\sum_{n=1}^N t^{(n)} x_i^{(n)}}{\sum_{n=1}^N t^{(n)}}}$$

$$\mu_{1i} = \frac{\text{number of examples where } x_i=1 \text{ and } t=1}{\text{number of examples where } t=1}$$

$$\mu_{2i} = \frac{\text{number of examples where } x_i=1 \text{ and } t=0}{\text{number of examples where } t=0}$$

Note: If the features are Categorical instead of binary, then the maximum likelihood solution for the parameters of the Categorical distribution would be the frequency of occurrence of that Category in the given class.

---

Example application: Text classification.

Given a document, which is a collection of words, classify the document into one of the 'k' topics.

ex: Politics, Sports, arts...

Assume there are M words in the Vocabulary.

Each document can be represented as a

$M$  dimensional feature vector

$$(x_1, \dots, x_M)$$

$$\begin{cases} 1 & \text{if first word is present in the document} \\ 0 & \text{otherwise.} \end{cases}$$

$$P(C_k | x) = P(x | C_k) P(C_k)$$

$$= \prod_{i=1}^M P(x_i | C_k) P(C_k)$$

Issue: Given limited number of examples per class, we might not see all words appearing at least once in all documents of every class.

if  $x_{1i}=0$  whenever  $t=1$ , then

$$M_{1i} = 0$$

This would make  $P(C_k | x)$  zero often.

How to avoid this?

Laplace Smoothing:-

$$M_{1i} = \frac{(\text{number of examples where } x_{1i}=1 \text{ and } t=1) + 1}{(\text{number of examples where } t=1) + 2}$$

→ If no example from a class, it reduces to a prior probability of  $y_2$ .



bias!

→ If all examples have  $x_i = 1$ , then

$$P(x_i = 0 | C_1) = \frac{1}{(\# \text{examples where } t=1) + 2}$$

→ the bias decreases, as we see the features set to 1 more often.

---

### Evaluation metrics for classification problem :-

#### 2-class problems :-

Positive / Negative

Accuracy =  $\frac{\text{number of correctly classified instances}}{\text{total number of instances}}$ .

Is this a good metric for all tasks?

Consider Cancer prediction. 9900 patients: Cancer = No  
100 patients: Cancer = Yes

classifier: default answer of "No"

↳ this classifier has 99% accuracy.

But this classifier is useless since it can't predict who has cancer.

Not every error has same cost!

Types of error this model can make:

- ① classify patients without cancer with "yes"
- ② classify patients with cancer as "no"

Second error is more dangerous than the first error.

Confusion matrix: True value

		Pos.	Neg.
Prediction	Pos.	True Positive	False Positive
	Neg.	False Negative	True Negative

In previous example, ① is False positive and  
② is false negative.

Example confusion matrix:

		100 examples	60+, 40 -
		true	
Prediction	Pos.	Pos.	Neg.
	Pos.	40	10
Neg.	Neg.	20	30
		60	40

$$\begin{pmatrix} 60 & 0 \\ 0 & 40 \end{pmatrix} \leftarrow \text{ideal confusion matrix.}$$

$$\begin{pmatrix} 30 & 30 \\ 40 & 30 \end{pmatrix} \leftarrow \text{bad confusion matrix.}$$

$$\text{Precision} = \frac{\text{true Pos}}{\text{true Pos} + \text{false pos}} = \frac{tp}{tp + fp}$$

$$\text{recall} = \frac{\text{true Pos}}{\text{true pos} + \text{false neg}} = \frac{tp}{tp + fn}$$

Precision  $\rightarrow$  Out of all positive predictions, how many of them are truly positive?

Recall  $\rightarrow$  Out of all positive examples in the data, how many of them are correctly identified as positive by the classifier?

\* Precision is important in face password systems.

\* Recall is important in email Spam classifier

↳ here good email  
↓  
Positive class.

It is easy to achieve high precision at the cost of low recall (and vice-versa).

① classifier which says default "yes"

↳ 100% recall, bad precision.

② classifier which says default "no"

↳ 100% precision, bad recall.

What if both precision and recall are important?

$$F_1 \text{ measure} = \frac{2 \cdot P \cdot R}{P + R}$$

P - Precision  
R - recall

Harmonic mean of precision and recall.

---

## You Should know!

1. Generalized linear models
2. logistic Sigmoid function
3. Softmax function
4. Gaussian discriminant analysis
5. independent and identically distributed (i.i.d)
6. Maximum likelihood approach
7. Quadratic discriminant analysis (QDA)
8. Naive Bayes assumption.
9. Gaussian Naive Bayes and Diagonal GDA .
10. Discrete Naive Bayes
11. Text classification
12. Laplace Smoothing.
13. Evaluation metrics for classification
  - Accuracy, confusion matrix, false positive, false negative, precision, recall,  $F_1$  measure.