

# Bayes ingenuo

Por:  
Ian Gabriel Cañas Fernández, 1092228

Profesor: Juan S. Pérez R.,  
Asignatura: INL367L, Secc 01

## Resumen:

En este laboratorio se estará implementando los conocimientos del clasificador de Naive Bayes aplicado a variables discretas, en las que se reconocerá la frecuencia de estos para así poder predecir si un pasajero a bordo al momento de la catástrofe del RMS Titanic sobrevivió o no al evento, para ello se considerará los datos obtenidos para cada pasajero, aplicando el teorema de Bayes para el acierto del clasificador.

## Ejercicios previos:

P5. Predecir la supervivencia.

En un dataset dado se solicita clasificar un set de datos de tripulantes del Titanic para reconocer si estos han sobrevivido o no, para ello, lo primero solicitado es su predicción mediante el uso de un clasificador Naive Bayes reconociendo que varias variables son discretas. La selección de las características asumió lo siguiente:

- EL ID del pasajero, su nombre, cabina y ticket son características únicas por pasajero.
- Pclass y Sex, son variables discretas representando la clase del boleto y el sexo de la persona, pues eran características representativas a la hora de seleccionar la supervivencia de las personas en contextos catastróficos. Incluyendo a la hora de seleccionar prioridades.
- La edad de la persona se ignorará, pues, aunque pudo haber prioridad por edades, no representa una distribución lo suficientemente simple al modelado.
- La frecuencia de los valores Embarked será ignorada porque la catástrofe aconteció luego de las embarcaciones a bordo, caso contrario no se presentaren en el listado.
- Se selecciona la cantidad de hermano y familiares, SibSp y Parch.

### E5.1 Naive Bayes; predicción de la supervivencia.

El primer método de machine learning utilizado para predecir la supervivencia de un grupo de personas se ha hecho una consideración de la frecuencia de personas que cumplen ciertas características y hayan resultado tanto muertas como vivas. Para no hacer un conteo para todas las posibles combinaciones, hacemos uso del teorema de bayes para más de una característica mediante la siguiente expresión:

$$\begin{aligned} P(C_K|x_1, x_2, \dots, x_d) &= P(C_K|\mathbf{x}) \\ &= \frac{P(C_K)P(\mathbf{x}|C_K)}{P(\mathbf{x})} \end{aligned}$$

Para el uso de las frecuencias se trabajó con los histogramas de cada característica según la clase que, en el caso, es la supervivencia (o no) de la persona. Calculamos las probabilidades a priori:

```
PVivo = sum(traindata.Survived==1)/length(traindata.Survived);  
PMuerto = sum(traindata.Survived==0)/length(traindata.Survived);
```

Los histogramas para cada clase individual

42	%% Survived, Pclass, SibSp, Parch.	
43	% sex,	
44		
45	liveMen = sum(vivos.SexN==0)/length(vivos.SexN);	
46	liveWomen = sum(vivos.SexN==1)/length(vivos.SexN);	
47		
48	deathMen = sum(muertos.SexN==0)/length(muertos.SexN);	
49	deathWomen = sum(muertos.SexN==1)/length(muertos.SexN);	
50		
51	[vivo_freqSex vivo_NSex] = hist(vivos.SexN, min(vivos.SexN):1:max(vivos.SexN));	
52	vivo_freqSex = vivo_freqSex./length(vivos.SexN);	
53		
54	[muerto_freqSex muerto_NSex] = hist(muertos.SexN, min(muertos.SexN):1:max(muertos.SexN));	
55	muerto_freqSex = muerto_freqSex./length(muertos.SexN);	
56		
57		
58	%% Histogramas para gente viva	
59		
60	[vivo_freqPclass vivo_NPclass] = hist(vivos.Pclass, min(vivos.Pclass):1:max(vivos.Pclass));	
61	[vivo_freqSibSp vivo_NSibSp] = hist(vivos.SibSp, min(vivos.SibSp):1:max(vivos.SibSp)*2);	
62	[vivo_freqParch vivo_NParch] = hist(vivos.Parch, min(vivos.Parch):1:max(vivos.Parch)*3);	
63		
64	vivo_freqPclass = vivo_freqPclass./length(vivos.Pclass);	
65	vivo_freqSibSp = vivo_freqSibSp./length(vivos.SibSp);	
66	vivo_freqParch = vivo_freqParch./length(vivos.Parch);	
67		
68	%% Histogramas para gente viva	
69		
70	[muerto_freqPclass muerto_NPclass] = hist(muertos.Pclass, min(muertos.Pclass):1:max(muertos.Pclass));	
71	[muerto_freqSibSp muerto_NSibSp] = hist(muertos.SibSp, min(muertos.SibSp):1:max(muertos.SibSp));	
72	[muerto_freqParch muerto_NParch] = hist(muertos.Parch, min(muertos.Parch):1:max(muertos.Parch)*3);	
73		
74	muerto_freqPclass = muerto_freqPclass./length(muertos.Pclass);	
75	muerto_freqSibSp = muerto_freqSibSp./length(muertos.SibSp);	
76	muerto_freqParch = muerto_freqParch./length(muertos.Parch);	
77		

Finalmente, conociendo las frecuencias, encontramos los likelihood de supervivencia para cada una de las características:

80	%% Likelihood por sexo	
81	minitest = testdata(:, :);	
82	ulu = mod(find(vivos.SexN(:)==minitest.SexN'), length(traindata.PassengerId))+1;	
83	% ulu = reshape(ulu, length(traindata.PassengerId), [])	
84		
85	i = vivo_NSex(:)==minitest.SexN';	103 %% Likelihood por SibSp
86	i = mod(find(i)-1, length(vivo_NSex))+1;	104
87	likeVivoSex = vivo_freqSex(i);	105 i = vivo_NSibSp(:)==minitest.SibSp';
88		106 i = mod(find(i)-1, length(vivo_NSibSp))+1;
89	i = muerto_NSex(:)==minitest.SexN';	107 likeVivoSibSp = vivo_freqSibSp(i);
90	i = mod(find(i)-1, length(muerto_NSex))+1;	108
91	likeMuertoSex = muerto_freqSex(i);	109 i = muerto_NSibSp(:)==minitest.SibSp';
92		110 i = mod(find(i)-1, length(muerto_NSibSp))+1;
93	%% Likelihood por Pclass	111 likeMuertoSibSp = muerto_freqSibSp(i);
94		112
95	i = vivo_NPclass(:)==minitest.Pclass';	113 %% Likelihood por Parch
96	i = mod(find(i)-1, length(vivo_NPclass))+1;	114
97	likeVivoPclass = vivo_freqPclass(i);	115 i = vivo_NParch(:)==minitest.Parch';
98		116 i = mod(find(i)-1, length(vivo_NParch))+1;
99	i = muerto_NPclass(:)==minitest.Pclass';	117 likeVivoParch = vivo_freqParch(i);
100	i = mod(find(i)-1, length(muerto_NPclass))+1;	118
101	likeMuertoPclass = muerto_freqPclass(i);	119 i = muerto_NParch(:)==minitest.Parch';
102		120 i = mod(find(i)-1, length(muerto_NParch))+1;
		121 likeMuertoParch = muerto_freqParch(i);
		122

Luego, aplicamos el teorema de bayes y comprobamos su eficiencia:

```

gVivo = PVivo.*likeVivoSex.*likeVivoPclass.*likeVivoSibSp.*likeVivoParch;
gMuerto = PMuerto.*likeMuertoSex.*likeMuertoPclass.*likeMuertoSibSp.*likeMuertoParch;

MLSurvived = nan(length(minitest.PassengerId),1);
minitest = [minitest table(MLSurvived)];

minitest.MLSurvived(:) = gVivo > gMuerto;
NaiveAccuracy = sum(minitest.MLSurvived==tester.Survived(:))/length(tester.Survived)

ind = tester.Survived(find(minitest.MLSurvived~=tester.Survived(:)));
falseNegative = sum(ind)/length(tester.Survived)
falsePositive = sum(~ind)/length(tester.Survived)

%%
% Si observamos los casos en que no se ha acertado el valor de
% supervivencia podemos observar que se presentan outliers, por lo que
% puede ser considerada una razón por la que tenemos fallos en los
% aciertos.

minitest(find(minitest.MLSurvived~=tester.Survived(:)), :)

```

## E5.2 Implementación en Matlab de kNN

De otra manera, se ha procurado hacer la predicción de supervivencia en base a las características, considerando la distancia euclideana de hacia las distintas clases. A continuación, se presenta el código en Matlab del algoritmo conocido implementado en Matlab.

```

20 %%%
21
22 MLSurvived = nan(length(testdata.PassengerId),1);
23 testdata = [testdata table(MLSurvived)];
24
25 for k = 1:2:7
26     for i = 1:length(testdata.PassengerId)
27         i;
28         d = ([testdata.SexN(i), testdata.Pclass(i), testdata.SibSp(i), testdata.Parch(i)] - [traindata.SexN(i), traindata.Pclass(i), traindata.SibSp(i), traindata.Parch(i)]).^2);
29         d = sqrt(d(:,1)+d(:,2)+d(:,3)+d(:,4));
30
31         [dis, pos] = mink(d, k);
32         testdata.MLSurvived(i) = mode(traindata.Survived(pos));
33     end
34
35 %%%
36 NaiveAccuracy = sum(testdata.MLSurvived==tester.Survived(:))/length(tester.Survived)
37 end

```

## R5.1 Resultados:

Para confirmar el rendimiento del algoritmo utilizaremos la métrica de tasa de aciertos de la predicción respecto a los datos obtenidos más tarde. Se presentan los resultados vistos por Naive Bayes y los vistos para kNN con diferentes cantidades de (k) vecinos cercanos.

```

NaiveAccuracy =

    0.9522

falseNegative =

    0.0144

falsePositive =

    0.0335

Elapsed time is 0.750710 seconds.

```

*Ilustración 1. Resultados de Naive Bayes.*

```

k =          k =          k =          k =
    1          3          5          7

kNNAccuracy = kNNAccuracy = kNNAccuracy = kNNAccuracy =
0.6914      0.5239      0.9139      0.9378

Elapsed time is 0.282120 seconds.

```

*Ilustración 2. Resultados para varias cantidades de vecinos mediante kNN.*

### A.5 Análisis:

En base a los resultados obtenidos mediante los dos métodos de clasificación, se pudo observar que el uso de Naive Bayes como clasificador es el más indicado en cuanto a características discretas se refiere, especialmente si estamos hablando de variables categóricas, pues este método es ciego al orden. Comparando los resultados se observa que, para mayores valores de k en kNN, mayor la eficiencia; sin embargo, en ningún caso kNN supera la eficiencia vista por Naive Bayes: un 95.22 %.