

# Regresión de mínimos cuadrados

Por:  
Ian Gabriel Cañas Fernández, 1092228

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

## Resumen:

En el presente laboratorio se planea computar los pesos de una regresión lineal usando el algoritmo Least Squares, a partir de el cual se hará una estimación de la salida del dataset obtenido. A partir del comportamiento del dataset que asumimos como lineal se desarrollará una regresión polinómica para observar la relación entre ambos métodos y las posibles variaciones de este último.

## Ejercicios previos:

### P6.1 Regresión lineal.

El objetivo principal de esta sección es el de computar los pesos de la regresión lineal a la simulación de un set de datos generado por el usuario. Este set sería segmentado en datos de entrenamiento y datos de prueba, siendo generado mediante la siguiente expresión:

$$r(t) = f(x) + e, \quad \text{donde } f(x) = w_1X + w_0$$

Se puede observar de antemano que la expresión simularía la medición ruidosa de un par de características con correlación lineal. Los valores seleccionados para la expresión fueron:

$$\begin{aligned} \mathbf{w} &= [-0.7 \ 220]; \\ X &\in \{1:1000\}; \\ e &\sim U(0,11). \end{aligned}$$

### P6.2 Predicción de las MPG.

Por otro lado, se espera la segmentación y predicción de un set de datos de Auto MPG Data Set. Se procurará computar los pesos de la regresión polinómica mediante el uso de Least Squares:

1. Considerando únicamente los dos *features* que arrojaron mejor desempeño.
2. Para 2 distintos sets de entrenamiento y pruebas.
3. Para 6 distintos grados de polinomios.

## E6.1 Naive Bayes; predicción de la supervivencia.

En la primera sección se genera, aleatoriza y segmenta el set de datos como se conoce de antemano mediante lo presentado a continuación:

```
4      %%
5      X = 1:1000;
6      X = X';
7      w = [-0.7 220];
8      Y = sum([ones(length(X), 1) X].*repmat(w, length(X), 1), 2);
9      Y = Y + normrnd(0, 11, length(X), 1);
10
11     data = [X Y];
12
13     dataset = data(randperm(length(X)), :); % "Dataset" porque son los datos con los que trabajaremos
14
15     i = round(length(X)*0.75);
16     D1 = ones(i, 1);
17
18     traindata = dataset(1:i, :);
19     testdata = dataset((i+1):end, :);
20
```

Más tarde, se genera una regresión lineal mediante la siguiente ecuación:

$$\mathbf{w} = (D^T D)^{-1} (D^T \mathbf{r})$$

Continuando con la aplicación de los pesos obtenidos al set de pruebas de los datos.

```
21 X_train = [ones(i,1) traindata(:,1)];
22 r = traindata(:, 2);
23
24 wA1 = inv(X_train'*X_train)*(X_train'*r)
25
26 X_test = [ones(length(testdata(:,1)),1) testdata(:,1)];
27 YML = X_test*wA1;
28
29 e1 = [YML testdata(:,2) (YML-testdata(:,2)).^2];
30 emean1 = mean(e1(:,3))
31
```

## E6.2 Naive Bayes; predicción de la supervivencia.

En la segunda sección, igualmente, empezamos aleatorizando los datos y seccionándolos en sets de entrenamiento y pruebas.

```
32 %%
33 data2 = readtable('auto-mpg.csv');
34 dataset2 = data2(randperm(length(data2.mpg)), :);
35 % "Dataset" porque son los datos con los que trabajaremos
36
37 i = round(length(dataset2.mpg)*0.75);
38
39 traindata2 = dataset2(1:i, :);
40 testdata2 = dataset2((i+1):end, :);
41
```

Preparamos las matrices de características junto a su vector bias para la regresión. Para la primera prueba, mediante el uso de regresión a tres dimensiones, se hace selección de las dos mejores características. El método utilizado para esta sección fue la de recorrer las combinaciones posibles de las permutaciones mediante forward search.

```
45 X2 = [ones(length(testdata2.cylinders(:)),1) testdata2.cylinders testdata2.displacement testdata2.weight testdata2.acceleration testdata2.model_year testdata2.origin];
46 X = [ones(1,1) traindata2.cylinders traindata2.displacement traindata2.weight traindata2.acceleration traindata2.model_year traindata2.origin];
47 r = traindata2.mpg;
48
49 sel = zeros(length(X(1,:)), 1);
50 l=2;
51 b=2;
52 a=0;
53 c=0;
54 while a==b
55     c = c + 1;
56     a=b;
57     sel = zeros(length(X(1,:)), 1);
58     for j = 2:length(X(1,:))
59         if a==j
60             continue
61         else
62             wA2 = inv(X(:, [1 a j]))'*X(:, [1 a j]))*X(:, [1 a j])'*r;
63
64             YML = X2(:, [1 a j])*wA2;
65
66             e2 = [YML testdata2.mpg (YML-testdata2.mpg).^2 1-abs(YML-testdata2.mpg)./length(testdata2.mpg)];
67             emean = mean(e2(:,3));
68             eprop = mean(e2(:,4));
69
70         end
71         sel(j) = eprop;
72     end
73     l = find(sel==max(sel));
74     l = min(l);
75     b=l;
76     if c==length(X(1,:))
77         break
78     end
79 end
```

Aplicamos regresión lineal y revisa error mínimo cuadrado:

```

85 wA2 = inv(X(:, [1 a b])'*X(:, [1 a b]))*X(:, [1 a b])'*r
86
87 YML = X2(:, [1 a b])*wA2;
88
89 e2 = [YML testdata2.mpg (YML-testdata2.mpg).^2 1-abs(YML-testdata2.mpg)./length(testdata2.mpg)];
90 emean2a = mean(e2(:,3))
91 eprop = mean(e2(:,4));

```

Finalmente, reprocesando los datos, se preparan las matrices de características y bias. Se implementa regresión de mínimos cuadrados para varios grados polinómicos entre hasta el grado seis (6):

```

93 %%
94
95 traindata2 = dataset2(1:i, :);
96 testdata2 = dataset2((i+1):end, :);
97
98 i = length(traindata2.cylinders);
99 j = length(testdata2.cylinders);
100
101 X2 = [ones(j,1) testdata2.cylinders testdata2.displacement testdata2.weight testdata2.acceleration t
102 X = [ones(i,1) traindata2.cylinders traindata2.displacement traindata2.weight traindata2.accelerati
103 r = traindata2.mpg;
104
105 for d = 2:length(X(1,:))
106     disp('Pesos y error para polinomio de grado ')
107     disp(d-1)
108
109     wA3 = inv(X(:,1:d)*X(:,1:d))*X(:,1:d)*r
110     mpgML = X2(:,1:d)*wA3;
111
112     e3 = [mpgML testdata2.mpg (mpgML-testdata2.mpg).^2];
113     emean3 = mean(e3(:,3))
114 end

```

## R6 Resultados:

Para la primera sección se obtuvo los siguientes resultados de pesos y error cuadrado promedio para varias corridas del código:

wA1 =	wA1 =	wA1 =	wA1 =
0.5444	-0.6428	-0.6717	0.2104
219.9985	220.0002	219.9995	219.9989
emean1 =	emean1 =	emean1 =	emean1 =
116.8407	94.5199	123.2305	128.3524

Para la segunda sección se trabaja con la mejor combinación de características para conocer su eficiencia:

```

wA2 =

    -14.9692
     -0.0069
      0.7773

emean2a =

    11.1919

```

Finalmente, se aplican las regresiones polinómicas:

Pesos y error para polinomio de grado 1	Pesos y error para polinomio de grado 2
---	---

wA3 =	wA3 =
43.4585	36.0394
-3.6466	-0.1168
	-0.0613

emean3 =	emean3 =
20.7290	20.9305

Pesos y error para polinomio de grado 3	Pesos y error para polinomio de grado 4
---	---

wA3 =	wA3 =
43.6048	39.6404
0.2233	0.1762
-0.0247	-0.0154
-0.0055	-0.0062
	0.2699

emean3 =	emean3 =
17.8851	18.5344

Pesos y error para polinomio de grado 5	Pesos y error para polinomio de grado 6
---	---

wA3 =

```
-16.8786
 0.1316
 0.0026
-0.0072
 0.1650
 0.7651
```

emean3 =

```
11.7478
```

wA3 =

```
-22.6243
 -0.0948
  0.0141
 -0.0072
  0.2400
  0.7801
 1.5054
```

emean3 =

```
11.2473
```

## A.5 Análisis:

En base a los resultados obtenidos mediante la regresión, se observó que los resultados de pesos obtenidos mediante la primera sección se veían directamente relacionados al error cuadrado promedio. Este error tiende más a parecer menor según los pesos de la predicción se asemejen a los computados anteriormente en la generación de los datos. Por lo que este método tiende a no ser muy alterado ante el ruido siempre y cuando este tenga su media en el cero.

Luego, comprobando lo obtenido con las dos mejores características y con varios grados de polinomios, se observó que los grados de polinomios afinan su eficiencia según se aumente dicho grado. Sin embargo, si se desea trabajar con pocas características representativas, es conveniente el uso de la mejor combinación, tomando en cuenta que trabajar con muchas características para este método no es computacionalmente tan pesado.