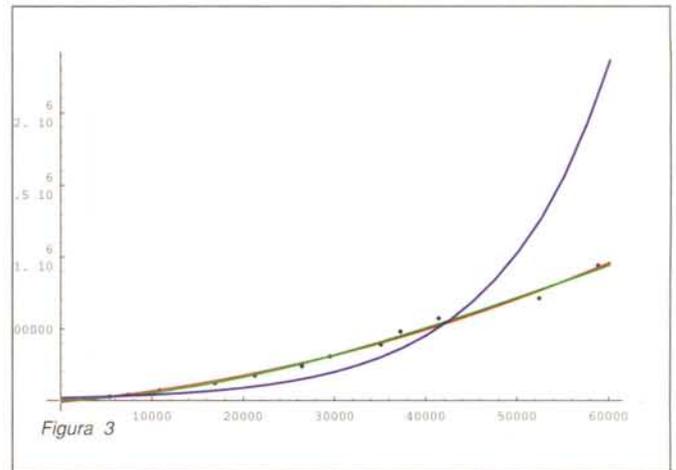


```
{... "prevedendo", "che", "le", "cinque",
"giare", "vecchie", "di", "coccio",
"smaltato", "che", "aveva", "in", "cantina",
"non", "sarebbero", "bastate", "a",
"contener", "tutto", "l", "olio", "della",
"nuova", "raccolta", "ne", "aveva",
"ordinata", "a", "tempo", "una", "sesta",
"pió", "capace", "a", "Santo", "Stefano" ...}
```

Correzione del testo

Studiando un poco le liste ottenute si evidenziano alcuni problemi: le lettere accentate sono tutte sbagliate a causa del diverso codice ASCII utilizzato, sono presenti sia maiuscole che minuscole, sono presenti caratteri non lettere (parentesi, numeri asterischi, etc). Vogliamo costruire un filtro che effettui le seguenti operazioni:

- corregge gli accenti
- trasforma le maiuscole in minuscole



• elimina i caratteri non lettere.
 Scriviamo dapprima la lista delle trasformazioni che correggono gli accenti. Questa lista è stata preparata a mano guardando un po' di testi, ci sono due correzioni diverse per la "ù" perché questa è rappresentata in due modi diversi in testi diversi.

```
In[1]:=
trasfdos={ "ó" -> "ù",
"é" -> "ù",
"ç" -> "è",
"ï" -> "ò",
"õ" -> "ï",
"ö" -> "à",
"ä" -> "é",
"ö" -> "à"};
```

Veniamo ora alle regole che trasformano le maiuscole in minuscole

```
In[2]:=
trasfM={
"A" -> "a", "B" -> "b", "C" -> "c", "D" -> "d",
"E" -> "e", "F" -> "f", "G" -> "g", "H" -> "h",
```

Un Modello matematico

Il fatto che il fitting migliore sia quello quadratico è un fenomeno curioso, di non facile spiegazione. In attesa di dedicarvi un'intera puntata, vi offro come antipasto un'idea di Giuseppe Fiorentino su cui stiamo lavorando.

Sia *u* il numero di tutte le parole della lingua in esame. Come prima approssimazione si assume che tutte le parole abbiano la stessa probabilità di presentarsi. Allora se ho trovato *t* parole dopo averne scandite *c* mi ne restano da trovare ancora *u-t* e la probabilità di trovarne di nuove non è 1 ma (*u-t*)/*u*. Si può scrivere allora l'equazione differenziale *dt/dc* = (*u-t*)/*u*.

Risolviamola con *Mathematica*

```
In[1]:=
DSolve[{t'[c]==1-t[c]/u,
t[0]==0}, t[c], c]
```

```
Out[1]=
      u
{{t[c] -> u - ——}
      Ec/u
```

```
In[2]:=
t[c_]:=Evaluate[t[c]/.%[1]]
```

E troviamo lo sviluppo in serie della funzione inversa *c(t)* che

è quella che ci interessa

```
In[3]:=
InverseSeries[Series[t[c], {c, 0, 5}]]
Out[3]=
      c2      c3      c4      c5
c + —— + —— + —— + —— + O[c6]
      2 u      3 u2      4 u3      5 u4
```

Si vede che il coefficiente del termine cubico è diviso per *u*² e quello del termine di quarto grado è diviso per *u*³. Nel nostro caso, con *u* dell'ordine del milione e *c* dell'ordine di qualche decina di migliaia significa che i termini di grado superiore al secondo possono essere trascurati.

Ovviamente l'ipotesi di distribuzione uniforme delle parole è molto imprecisa; possibili sviluppi di ricerca consistono nel trovare modelli più realistici e nel confrontarne le previsioni con i risultati sperimentali ottenuti su un corpus che sia una frazione consistente di una lingua. Con le lingue moderne questo è di difficile realizzazione ma per il Latino e il Greco ci stiamo organizzando...

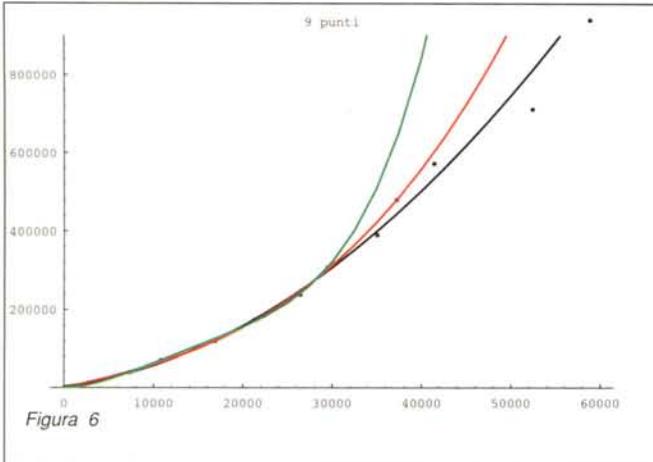


Figura 6

09	SENILITA	{69043,	47.9,	307431,	29520}
10	C1	{82749,	90.9,	390180,	35085}
11	MALAV	{90464,	44.5,	480644,	37249}
12	C5	{92495,	88.4,	573139,	41429}
13	C6	{139301,	143.5,	712440,	52432}
14	PROMESSI	{228881,	121.9,	941321,	58834}

Si nota che i tempi di elaborazione sono in generale accettabili (2 minuti per elaborare i "Promessi Sposi").

Ovviamente un programma C ottimizzato richiederebbe una frazione di questo tempo di elaborazione ma un tempo uomo almeno dieci volte maggiore per la programmazione.

Fitting

Le ultime due colonne mostrano come variano rispettivamente c la cardinalità dell'insieme delle parole scandite e t la cardinalità dell'insieme delle parole trovate.

Il plottaggio di $p1$ dà un'idea di come cresce la funzione $c(t)$.

```
In[1]:=
listplot=ListPlot[p1,
PlotRange->All];
```

(Figura 2)

Ovviamente la crescita non è lineare. Per avere un'idea più precisa di come cresce la funzione $c(t)$ proviamo a fare qualche approssimazione ai minimi quadrati. Iniziamo con un polinomio di secondo grado:

```
In[2]:=
f2=Fit[p1, {1,x,x^2}, x]
Out[2]=
-10518.2 + 5.58235 x + 0.000176063 x^2
```

e di terzo grado:

```
In[3]:=
f3=Fit[p1, {1,x,x^2,x^3}, x]
Out[3]=
0.000360445 + 2.46747 t +
0.000318208 t^2 - 1.62262 10^-9 t^3
```

Un altro tentativo può essere la funzione esponenziale $a e^{bx}$. Per fare questa approssimazione è necessario applicare l'identità $\text{Log}(a e^{bx}) = \text{Log } a + b x$ e quindi ricavare $\text{Log } a$ e b da un fit lineare del logaritmo naturale dei dati.

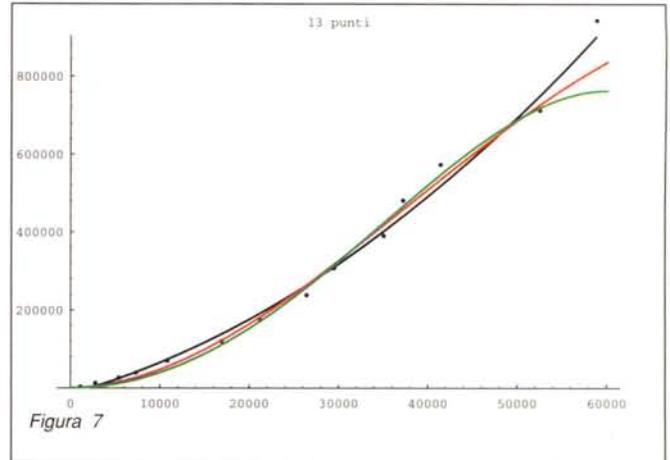


Figura 7

```
In[4]:=
lp1=
Map[{{#[[1]],Log[N#[[2]]]}&,p1];
In[5]:=
fe=Expand[Exp[Fit[lp1, {1,x}, x]]]
Out[5]=
16562.8 E^0.0000826762 t
Disegniamo insieme i tre grafici e i punti sperimentali:
In[6]:=
plot=Plot[{f2,f3,fe},{x,0,60000},
PlotStyle->{Red,Green,Blue}];
Show[listplot,plot];
```

(Figura 3)

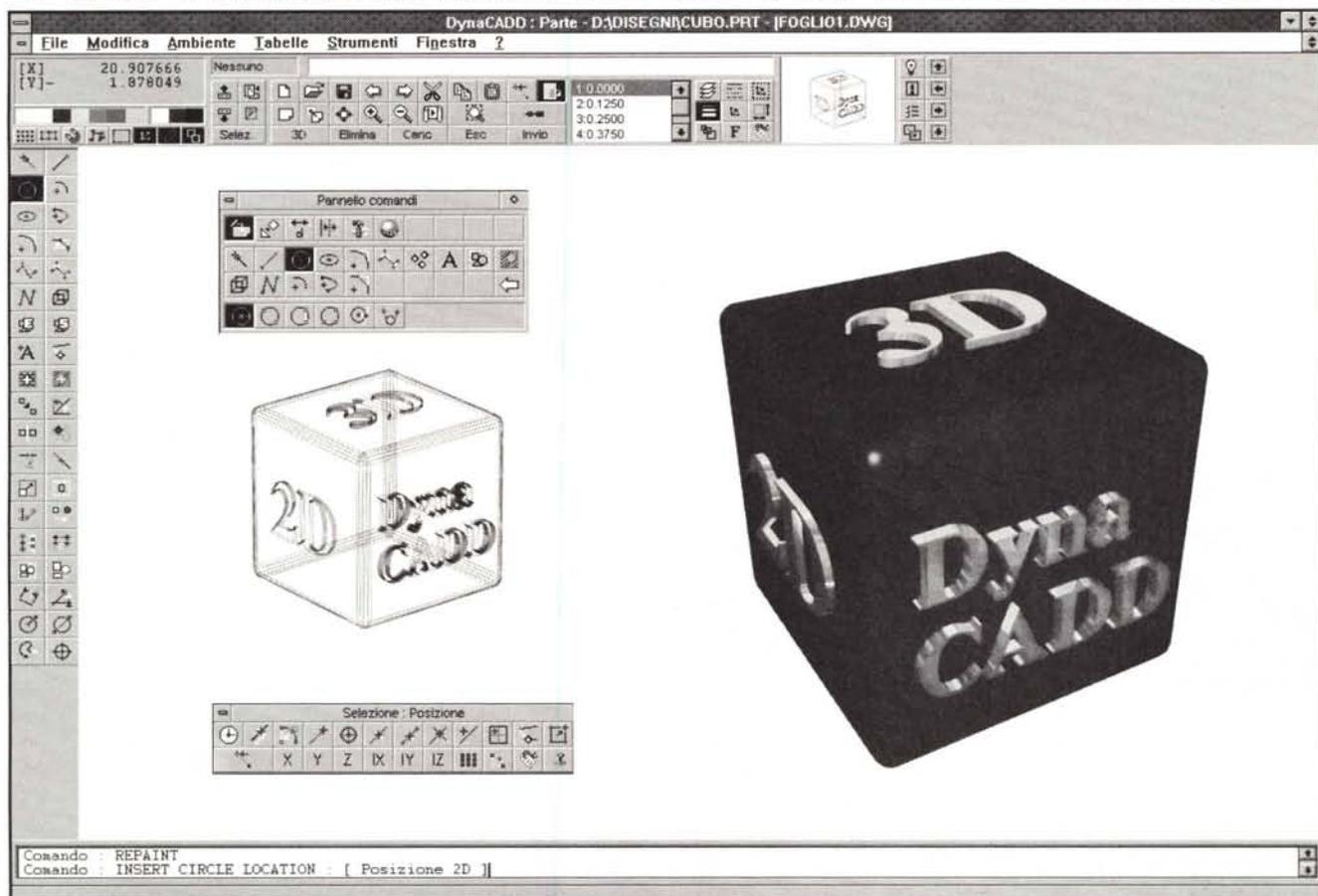
Si nota che il fitting esponenziale sballa completamente mentre i due fitting polinomiali sono abbastanza accurati. Per decidere quale sia il grado più appropriato per il polinomio proviamo ad approssimare solo i primi punti e vediamo come vengono previsti gli altri. Il seguente programma approssima con polinomi di secondo, terzo e quarto grado, partendo con 6 punti e arrivando fino a 13.

```
In[7]:=
Do[
f2=Fit[Take[p1,np1],
{1,x,x^2},x];
f3=Fit[Take[p1,np1],
{1,x,x^2,x^3},x];
f4=Fit[Take[p1,np1],
{1,x,x^2,x^3,x^4},x];
plot=Plot[{f2,f3,f4},{x,0,60000},
PlotStyle->{Black,Red,Green},
DisplayFunction->Identity];
Show[listplot,plot,
DisplayFunction->
$DisplayFunction,
PlotRange->{0,900000},
PlotLabel->
ToString[np1]<>"punti",
{np1,6,13}];
```

Le figure da 4 a 7 mostrano i risultati per $np1=6,7,9,13$, rispettivamente. È evidente che l'unica approssimazione che permette di fare previsioni ragionevoli è quella quadratica.

DynaCADD®

I CAD PROFESSIONALI 2D-3D PER DOS E WINDOWS



DynaDesigner

WIN

CAD 2D per Windows - Associativo - 256 layer - 24 tipi di primitive - Oltre 2000 comandi tramite interfaccia grafica o linea comandi - Undo e Redo infiniti - DXF in/out - Autorecover in grado di recuperare tutto il lavoro in caso di crash - Cursore intelligente con 8 possibilità di snap - Programmabile in C tramite il sistema di sviluppo e conversione dei font True Type tramite il Font Editor (non inclusi).

Lit. 550.000

DynaCADD

WIN

CAD 2D/3D per Windows - Le funzioni di DynaDesigner e in più:
 • Funzioni per la creazione di solidi • Oltre 2500 comandi • Viste tridimensionali illimitate • Rendering a 24 bit con shading, shadowing (con algoritmi di Pixar) e texture mapping • Programma per l'editing bitmap incluso • Servizio di Hotline gratuito.

Lit. 1.200.000

DynaCADD 2D-3D

DOS

Cad 2D e 3D wireframe per DOS - 256 layer - 13 tipi di primitive incluse curve di Bézier e B-splines - Uscita su stampanti, plotter e dispositivi Postscript - Font vettoriali ed editor di font inclusi - Interscambio file DXF 2D e 3D sia in lettura che in scrittura - Help in linea - Viste tridimensionali multiple - Quotature automatiche - Precisione a 16 cifre - Interfaccia utente semplice ed intuitiva - Servizio di Hotline gratuito.

Lit. 250.000



IMPORTATORE E DISTRIBUTORE ESCLUSIVO PER L'ITALIA: STUDIO NUOVE FORME SRL

Via Mancinelli, 19 - 20131 Milano - Tel. 02/26143833 r.a. - Fax 02/26147440 - DynaCADD Hot-line 02/26149649