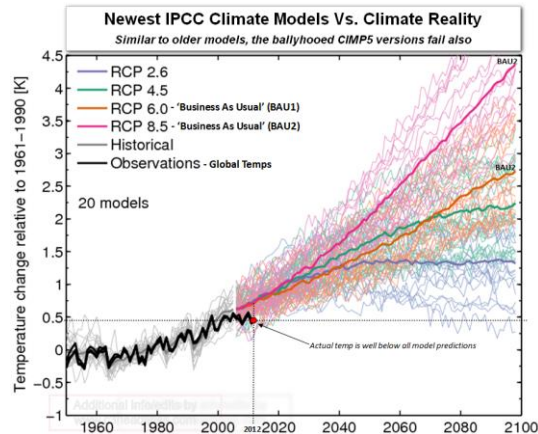


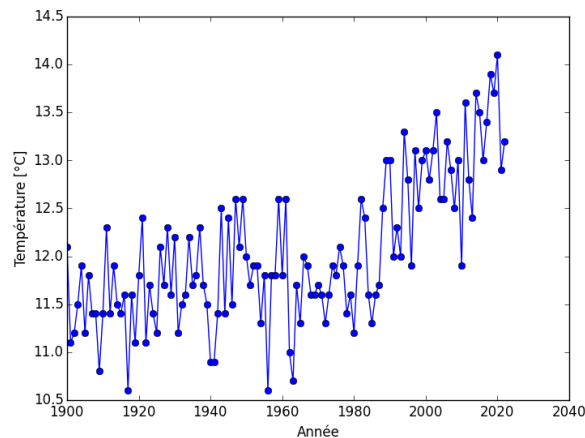
Modèle de prédiction de la température annuelle moyenne en France



Préparation des données

On s'intéresse à l'évolution de la moyenne de la température annuelle en France depuis l'année 1900, première année pour laquelle des relevés réguliers ont été réalisés.

Le fichier *Temp_annuelle_France.xlsx* comprend deux colonnes dans lesquelles on retrouve l'évolution de la température moyenne annuelle en fonction de l'année. Ces deux jeux de données seront appelés par la suite *Annee* et *Temp*. Les données sont tracées sur le graphique suivant :



Q1 : Après avoir calculé la moyenne et l'écart type des données *Annee* et *Temp*. Centrer ces données en stockant pour chacune d'entre elles la moyenne et écart type. On rappelle que le centrage des données est donné par la formule suivante : $X = \frac{X - \text{moyenne}}{\text{ecart type}}$.

Une régression linéaire multiple est envisagée, avec un modèle polynomial. Plusieurs degrés seront comparés (de 0 à 6) afin de valider l'évolution de la température, selon plusieurs modèles.

Régression avec Régression Linéaire Multiple

Afin de pouvoir effectuer l'apprentissage sur les données étudiées ici, nous allons devoir programmer trois fonctions intermédiaires :

- Fonction modèle h prenant en arguments une matrice X et le vecteur de paramètres W et qui retourne la valeur : $h(\text{Annee}, W)$, où h est la fonction polynomiale.
- Fonction coût J prenant en arguments la matrice des entrées, le vecteur des sorties et le vecteur des paramètres et calculant le coût (ou l'erreur) de la fonction modèle.
- Fonction *Gradient*, prenant les mêmes arguments que la fonction précédente et qui sera utilisée dans l'algorithme de descente de Gradient. Le calcul du gradient est donné par la formule : $\text{Gradient}(X, Y) = \frac{1}{N} \cdot X^T \cdot (\tilde{Y} - Y)$ avec $N = \text{Nombre d'éléments de } X$ et \tilde{Y} la fonction approchée de Y à l'aide du modèle polynomial.

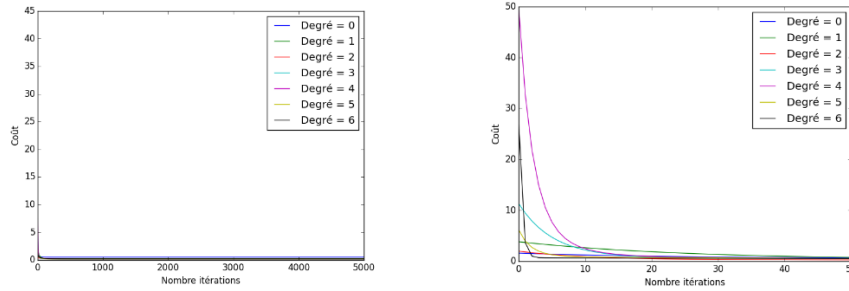
Q2 : Coder ces trois fonctions, ainsi qu'une fonction globale *Regression_Lineaire* prenant en arguments la matrice *Annee* des observations, le vecteur *Temp* des étiquettes, le degré d du modèle du modèle polynomial, le nombre d'itérations maximales, le taux d'apprentissage, et qui retourne le vecteur W des paramètres optimisant le coût, ainsi que l'historique de ce coût à chaque passage de boucle sous forme de liste.

Remarque : La méthode de la descente du gradient utilise la méthode suivante dans la boucle d'itération : $W = W - \alpha \cdot \text{Gradient}(X, Y, W)$ avec α le taux d'apprentissage.

Remarque : On prendra garde à ne pas oublier dans la fonction *Regression_Lineaire* l'ajout d'une colonne de 1 dans la matrice *Annee* pour prendre en compte le biais dans le vecteur W , ainsi qu'un nombre de colonnes correspondant au degré du modèle polynomial choisi (par exemple, on aura une colonne supplémentaire contenant les carrés des entrées pour un degré 2, etc.). Cette fonction est donnée par le code suivant :

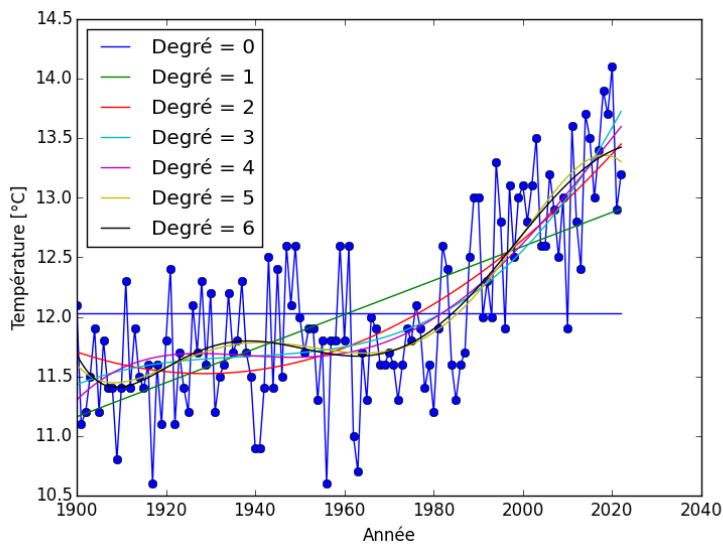
```
def Degre_multiple(X, d):
    N = np.shape(X)[0]
    U = np.ones((N, 1))
    if d==0:
        R=U
    else:
        R=np.hstack((U,X)) # np.hstack Permet d'empiler des séquences horizontalement
        k=2
        while k<=d:
            R=np.hstack((R,X**k))
            k+=1
    return R
```

En utilisant les fonctions précédemment créées, on trace l'évolution du coût (qui est, on le rappelle pour une régression linéaire, une erreur quadratique) en fonction du nombre d'itérations. La convergence est estimée atteinte lorsque ce coût tend vers une asymptote horizontale. La courbe obtenue pour 5000 itérations est obtenue ci-après.



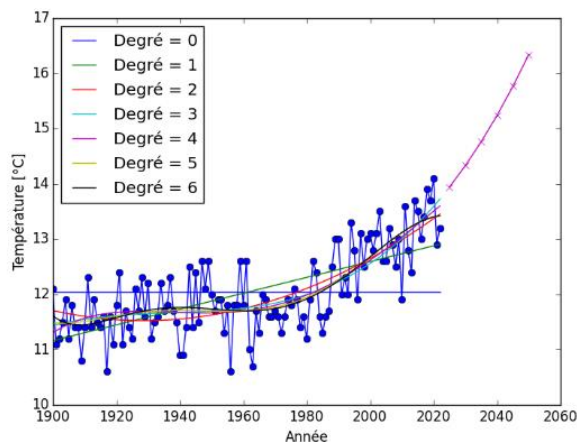
Q3 : Coder une fonction *Prediction* qui prend en entrées une matrice X et le vecteur de paramètres W et qui retourne la fonction modèle h avec en arguments une matrice X (avec la bonne taille en fonction du degré) et le vecteur de paramètres W .

Les différentes prédictions sont visibles sur la courbe ci-dessous.



Q4 : Déterminer le degré du modèle polynomial qui vous semble le mieux adapté.

Q5 : Calculer, avec le modèle choisi, les prédictions pour les années 2025, 2030, 2035, 2040, 2045 et 2050.



Utilisation de la bibliothèque *scikit-learn*

Il est possible d'utiliser des fonctions toutes faites, au lieu de re-coder l'ensemble des fonctions. En effet, les modèles linéaires sont implémentés dans le module *linear_model*. La régression linéaire elle-même est implémentée dans la classe *linear_model.LinearRegression*.

Pour la régression linéaire polynomiale, il est nécessaire de préparer les données en générant des données d'entrées compatible avec une régression polynomiale grâce à *PolynomialFeatures*, utilisé ainsi :

```
from sklearn.preprocessing import PolynomialFeatures
transformation_polynomiale = PolynomialFeatures(degree = degre)
X_t = transformation_polynomiale.fit_transform(X)
```

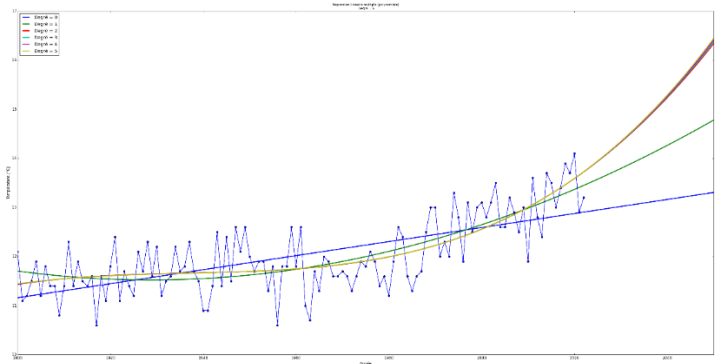
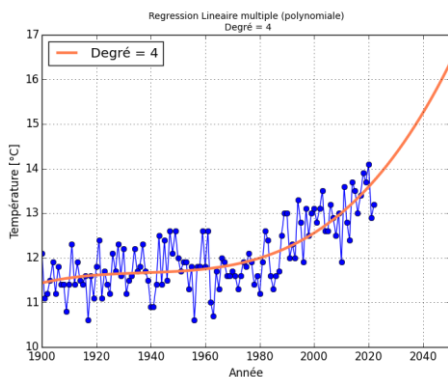
Dans les lignes de code ci-dessus, *degre* est un entier indiquant le degré du polynôme, tandis que *X_t* sera le résultat de la transformation générant une matrice incluant un nombre de colonne correspondant au degré. *X* est le vecteur des entrées.

Pour la régression polynomiale, on instancie un objet de la classe de modèle qui nous intéresse, ici *LinearRegression()*. De même, on entraînera cet objet sur les données d'entraînement (*(X_t, Y)*) et non (*X, Y*) avec la méthode *.fit*.

Enfin, pour prédire des résultats, une fois le modèle entraîné, on utilisera la méthode *.predict*. On remarque alors que quel que soit le type d'apprentissage utilisé, les méthodes pour entraîner le modèle puis faire des prédictions restent les mêmes...

Enfin, il est possible d'évaluer les performances du modèle d'apprentissage grâce au module *metrics*, qui s'importe. On utilisera la fonction *mean_squared_error* qui prend en arguments *Y*, les étiquettes à prédire, et *Ym*, les sorties prédites par le modèle. On aura alors l'erreur quadratique entre les données prédites et à prédire par le modèle entraîné.

Q6 : Compléter les étapes 2, 3, 4 et 6 afin de permettre d'entraîner le modèle, de prédire l'évolution de la température en fonction des années de 1900 à 2050 et d'évaluer les performances du modèle d'apprentissage en calculant l'erreur quadratique. On pourra utiliser le meilleur degré du polynôme choisi (0 à 6) déterminé dans la partie précédente.



```

1 # -*- coding: utf-8 -*-
2 import os
3 os.chdir('CHEMIN DES FICHIERS')
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 # Extraction des données et création des 2 jeux de données
8 dataset=np.loadtxt('Temp_annuelle_France.csv', skiprows=1,delimiter=',')
9 Annee=np.array([dataset[i][0] for i in range(len(dataset))])
10 Annee = Annee.reshape((len(Annee)),1)
11 Temp=np.array([dataset[i][1] for i in range(len(dataset))])
12 Temp = Temp.reshape((len(Temp)),1)
13
14 # TRACE Dataset
15 plt.figure(1)
16 plt.plot(Annee,Temp,'bo-')
17 plt.xlabel('Année')
18 plt.ylabel('Température [°C]')
19
20 # Calcul des moyennes et écart types
21 mu_Annee =
22 sigma_Annee =
23 mu_Temp =
24 sigma_Temp =
25
26 # Centrage des données
27 Annee =
28 Temp =
29
30 def h(X,W):
31     return
32
33 def J(X,Y,W):
34     N=np.shape(X)[0]
35     cout =
36     return
37
38 def Gradient(X,Y,W):
39     N = np.shape(X)[0]
40     return
41
42 def Degre_multiple(X,d):
43     N = np.shape(X)[0]
44     U = np.ones((N,1))
45     if d==0:
46         R=U
47     else:
48         R=np.hstack((U,X)) # np.hstack Permet d'empiler des séquences horizontalement
49         k=2
50         while k<=d:
51             R=np.hstack((R,X**k))
52             k+=1
53     return R
54
55 def Regression_Lineaire(X,Y,degre,nbre_iter=1000,alpha=0.01):
56     W = np.random.randn(degre+1,1)
57     X = Degre_multiple(X,degre)
58     Histoire_cout = [J(X,Y,W)]
59     k=0
60     while k<nbre_iter :
61
62         k+=1
63     return W,Histoire_cout
64
65
66 L_W = []
67 L_cout = []
68
69 for k in range(0,7):
70     W,ch = Regression_Lineaire(Annee,Temp,k,nbre_iter=5000,alpha=0.02)
71     L_W.append(W)
72     L_cout.append(ch)
73
74 plt.figure(2)
75 for k in range(len(L_cout)):
76     labell = 'Degré = {}'.format(k)
77     plt.plot(L_cout[k],label=str(labell))
78 plt.xlabel("Nombre itérations")
79 plt.ylabel("Coût")
80 plt.legend()
81
82 plt.figure(3)
83 for k in range(len(L_cout)):
84     labell = 'Degré = {}'.format(k)
85     plt.plot(L_cout[k],label=str(labell))
86 plt.xlabel("Nombre itérations")
87 plt.ylabel("Coût")
88 plt.xlim(0,50)
89 plt.legend()
90
91 def Prediction(X,W):
92     degre = len(W)-1
93
94     return
95
96 plt.figure(1)
97 for k in range(len(L_W)):
98     ym = Prediction(Annee,L_W[k])
99     X = Annee*sigma_Annee + mu_Annee
100     ym = ym*sigma_Temp + mu_Temp
101     labell = 'Degré = {}'.format(k)
102     plt.plot(X,ym,label=str(labell))
103 plt.legend(loc="upper left")
104
105 # Un degré 4 semble être la solution offrant le modèle le plus proche
106 # prédiction températures en 2030, 2050 et 2100...
107
108 Xp = np.array([[2025,2030,2035,2040,2045,2050]])
109 Xpn = # On normalise pour coller au modèle
110 yp = # On prend le modèle de degré ??
111 yp =
112
113 plt.figure(1)
114 plt.plot(Xp,yp,'mx-')

```

```

1 # -*- coding: utf-8 -*-
2 import os
3 os.chdir('CHEMIN DES FICHIERS')
4 from sklearn.linear_model import LinearRegression
5 from sklearn.preprocessing import PolynomialFeatures
6 from sklearn.metrics import mean_squared_error
7 import numpy as np
8 import matplotlib.pyplot as plt
9
10 # Etape 1 : Recuperation du dataset
11 dataset=np.loadtxt('Temp_annuelle_France.csv', skiprows=1,delimiter=';')
12 Annee=np.array([dataset[i][0] for i in range(len(dataset))])
13 X = Annee.reshape((len(Annee)),1)
14 Temp=np.array([dataset[i][1] for i in range(len(dataset))])
15 Y = Temp.reshape((len(Temp)),1)
16
17 plt.figure(1)
18 plt.plot(X,Y,'bo-')
19 plt.xlabel('Année')
20 plt.ylabel('Température [°C]')
21
22 # Etape 2 : Preparation des donnees
23 degre =
24
25
26
27 # Etape 3: Definition et entraînement du modèle
28 Modele_RP =
29
30
31 # Etape 4 : Prediction
32 x_new_min = 1900
33 x_new_max = 2050
34 Xp = np.linspace(x_new_min, x_new_max, 2050-1900)
35 Xp = Xp[:,np.newaxis]
36 Xp_t =
37 Yp =
38
39 # Etape 5 : Graphiques
40 plt.figure(1)
41 plt.plot(Xp, Yp,label=str(label), color='coral', linewidth=3)
42 plt.grid()
43 label = 'Degré = {}'.format(degre)
44 plt.xlim(x_new_min,x_new_max)
45 title = 'Degré = {}'.format(degre)
46 plt.title("Regression Lineaire multiple (polynomiale) \n " + title, fontsize=10)
47 plt.xlabel('Année')
48 plt.ylabel('Température [°C]')
49 plt.legend(loc="upper left")
50 plt.show()
51
52 # Etape 6 : Calculs des erreurs
53 Ym =
54 rmse =
55 print('RMSE: ', rmse)

```