

Méthodes Statistiques pour la Finance

Victor MOINE



Contents

I	Analyse en composantes principales.	2
1	Présentation du problème.	3
2	Données du problème.	3
3	Résolution.	4
4	Interprétation.	4
4.1	Représentation des variables sur l'axe 1 :	4
4.2	Projection des individus sur l'axe 1 :	5
4.3	Représentation des variables dans le plan axe 1 / axe 2 :	6
4.4	Projection des individus dans le plan axe 1 / axe 2 :	7
II	Régression linéaire : régression MCO et Ridge.	8
5	MCO : présentation du problème & hypothèses	8

6	Données du problème.	9
7	Résolution et interprétation.	10
7.1	Intervalles de confiance sur θ_j .	11
7.2	Sélection de variables.	12
8	Régression Ridge.	13
8.1	Expression de l'estimateur Ridge.	13
8.2	Comment fixer λ : la cross-validation.	14
III	Portefeuille de Markowitz	16
9	Présentation du problème.	16
10	Données du problème.	16
11	Cas 1 : en présence d'actif sans risque.	16
11.1	Mise en équation.	16
11.2	Solution.	17
12	Cas 2 : sans ASR.	18
12.1	Mise en équation.	18
12.2	Solution.	18
12.3	Résultat.	18
12.3.1	Constatation.	18
12.3.2	Exemple appliqué du trader.	19
IV	Annexe : code des 3 parties (ACP - Régression MCO - Régression Ridge - Markowitz) - cf pages suivantes	20

Abstract

13 pages (en enlevant les pages de garde, les espaces dûs à l'inclusion de graphiques sur Latex) / Les liens des données sont en commentaires des scripts Matlab

Part I

Analyse en composantes principales.

1 Présentation du problème.

Le but de l'analyse en composantes principales est de représenter de manière parcimonieuse et pertinente un nuage de points vivant dans un espace de grande dimension. C'est une méthode de projection linéaire pour diminuer le nombre de paramètres. En effet, l'ACP va projeter les données sur un espace de dimension moins élevée par rapport à l'espace d'origine.

2 Données du problème.

On dispose d'un tableau de n lignes et p colonnes.

- Une ligne i : (x_i^1, \dots, x_i^p) = individu i
- Une colonne j : (x_1^j, \dots, x_n^j) = individu j

Dans notre exemple:

- les individus seront les jours de trading: du 02/11/2015 au 22/01/2016
- les variables seront les rendements des actifs : 36 actifs du CAC40 (et non 40, certaines données par exemple pour Alstom n'étant pas disponible sur toute la période)

On va considérer des versions centrées et réduites des variables: on va effectuer les opérations suivantes sur les variables :

$$\bar{x}^j = \frac{1}{n} \sum_{i=1}^n x_i^j$$
$$x_i^j \leftarrow \frac{x_i^j - \bar{x}^j}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i^j - \bar{x}^j)^2}}$$

En effet ici, bien que toutes les données soient "homogènes à des rendements" (ie: elles sont à peu près de même grandeur), on peut dans certain cas plus généraux réaliser une ACP sur des variables hétérogènes. Normer les variables, permet ainsi de réduire l'effet de l'ordre de grandeur de variables de natures différentes.

3 Résolution.

Comme dans le cours, les 3 étapes d'une ACP sont :

- Etape 1 : construction de la matrice de corrélation empirique entre les variables V
- Etape 2 : recherche des plus grandes valeurs propres de V et des vecteurs propres associés
- Etape 3 : représentation par projection et interprétation

Algorithm 1 extrait de acp.m - Suite : cf annexe

```

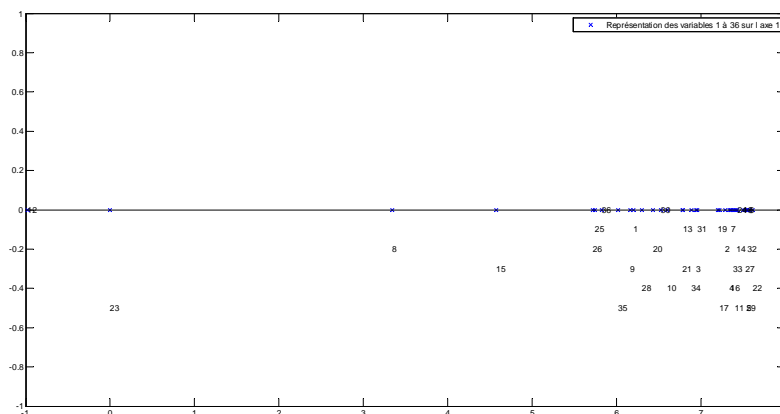
1 % Centrage et normalisation des points variables
2 for j=1:p
3     mu = (1/n)*sum(X(:,j)); % moyenne du point variable j
4     Y = X(:,j)-mu;
5     va = (1/n)*(Y')*Y;
6     X(:,j) = (X(:,j)-mu)/sqrt(va);
7 end

```

4 Interprétation.

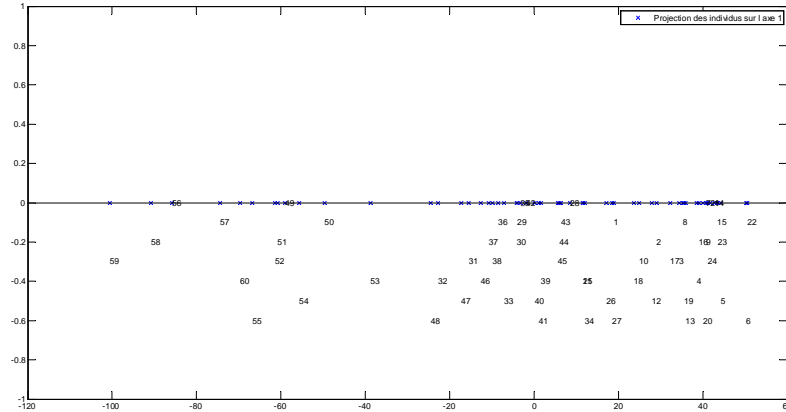
4.1 Représentation des variables sur l'axe 1 :

Les variables sont numérotées de 1 à 36 (le numéro est placé en dessous ou au dessus de la projection de la variable sur l'axe 1).



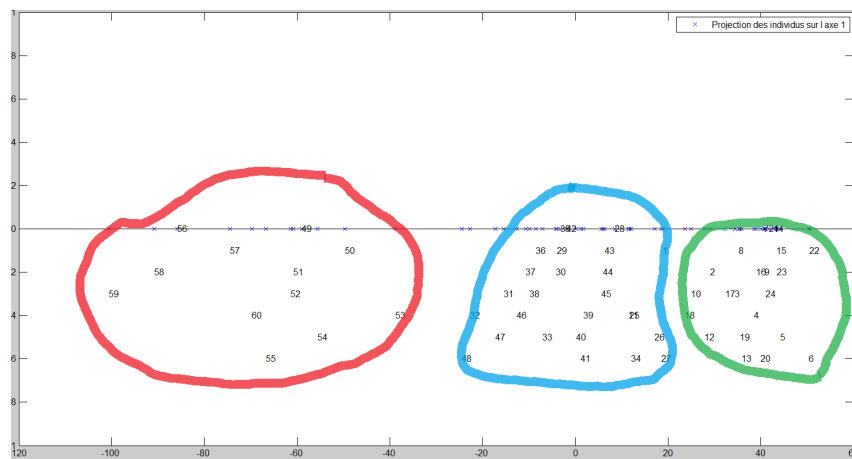
L'axe 1 est le vecteur propre $v_1 = [-0.0577, 0.0638, -0.0527, \dots]$ qui est bien de la forme $[\alpha + / - \epsilon, \alpha + / - \epsilon, \dots]$ d'après Matlab.

4.2 Projection des individus sur l'axe 1 :



On remarque, comme le confirme la théorie, que l'axe 1 est la meilleure représentation en dimension 1 du nuage des individus, au sens où il maximise la dispersion, car les projections des “jours de trading” sont en effet bien étalés sur l'axe 1.

On remarque par exemple que puisque la composante Accor (variable 1) est positive sur l'axe 1, alors pour tous les jours qui se projettent positivement, la corrélation entre ces jours et Accor est positive. Ou réciproquement, pour tous les jours tels que leur corrélation avec Accor est positive, alors s'ils se projettent négativement sur l'axe 1 cela implique que Accor était ≤ 0 ces jours-là.

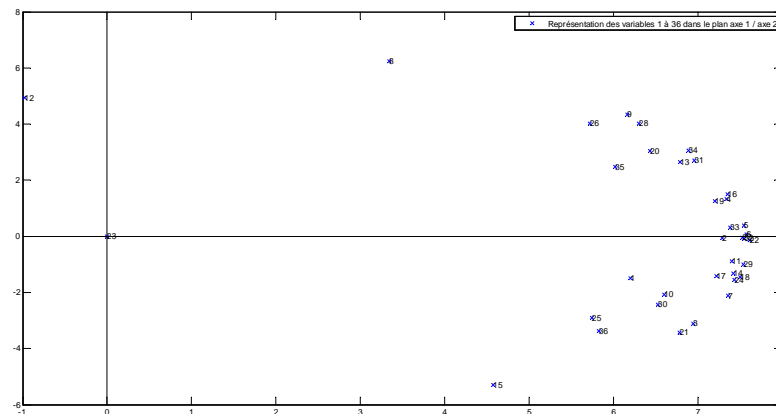


On peut aussi distinguer 3 groupes : les jours de trading (numérotés de 1 à 60 à partir de la date initiale du 02/11/2015) où le marché est

- bull : pour les jours entourés en vert (se projettent très positivement sur l'axe 1)
- bear : pour les jours entourés en rouge
- flat : pour les jours entourés en bleu

On constate aussi un effet de clustering : des jours consécutifs appartiennent au même groupe.

4.3 Représentation des variables dans le plan axe 1 / axe 2 :



On constate que l'axe 2 oppose 2 types de variables: les variables

- 8/26/35/9/28/20/13/34/31/19/4/12
- 25/36/1/30/10/21/7/17

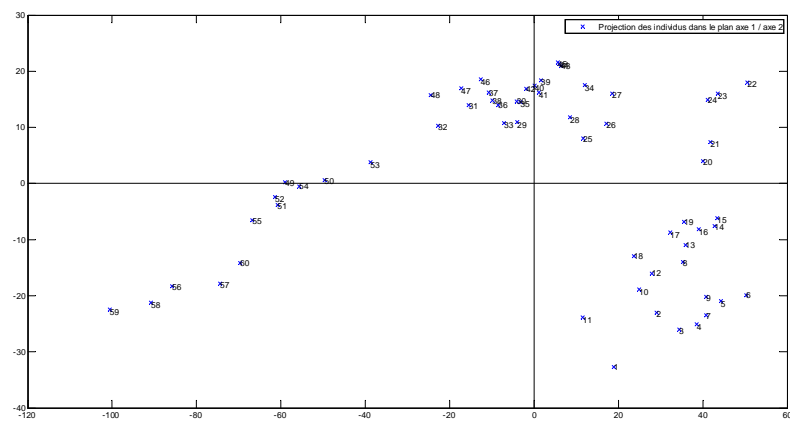
correspondant respectivement aux actifs :

- CapGemini/Publicis/Veolia/Axa/Renault/Legrand/Engie/Peugeot/St Gobain/Klepierre/Airbus/Bouygues
- Orange/Vivendi/Accor/Sanofi/Vinci/LVMH/Carrefour/Kering

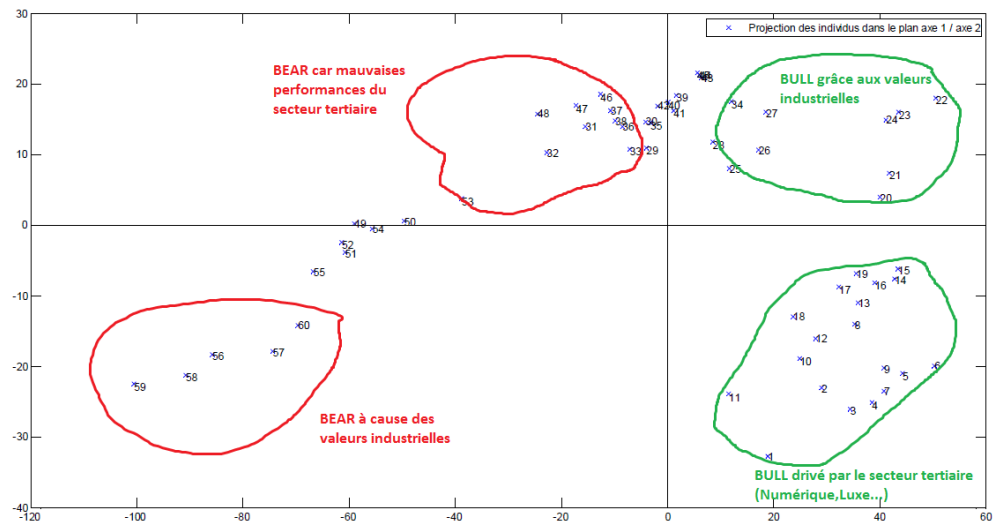
On constate donc que l'axe 2 distingue globalement les valeurs industrielles (Veolia, Renault, Legrand, Peugeot, St Gobain, Airbus, Bouygues) aux valeurs du secteur tertiaire des services (du numérique : Orange, Vivendi / du luxe : LVMH, Kering): l'axe 2 est un axe sectoriel.

4.4 Projection des individus dans le plan axe 1 / axe 2 :

Voici la meilleure représentation de notre nuage de points en dimension 2 : là encore on remarque un effet de clustering (les jours 56,57,58,59,60 sont proches, de même que 46,47,48,...) ce qui est logique.



On peut observer ceci :



En effet, par exemple pour les valeurs du secteur tertiaire telles que LVMH (variable 21), celles-ci sont très négatives sur l'axe 2, donc la corrélation entre les jours 4/5/6/7 et ces valeurs du secteur tertiaire sont négatives: les jours 4/5/6/7 (correspondant aux 6,7,8,9 novembre 2015) sont donc des journées BULL drivées par les valeurs du secteur tertiaire (à 1 jour qui se projette négativement correspond une valeur du secteur tertiaire positive car la corrélation est négative).

En résumé, il y a 4 types de journées: BULL ou BEAR et ceci A CAUSE(ou GRACE) aux valeurs INDUSTRIELLES ou du SECTEUR TERTIAIRE.

Part II

Régression linéaire : régression MCO et Ridge.

5 MCO : présentation du problème & hypothèses

La régression linéaire a pour but d'expliquer une variable Y à partir d'une combinaison linéaire de variables explicatives X , de sorte que

$$Y = \phi\theta + \varepsilon$$

où :

- Y est le vecteur des n observations
- ϕ est la matrice (n, p) des régresseurs
- θ est le vecteur inconnu de dimension p
- ε est un bruit

On fait les 2 hypothèses suivantes :

- H1 : $\text{rang}(X) = p$
- H2 : $\varepsilon \sim N(0, \sigma^2 I_p)$

L'hypothèse H2 n'est pas nécessaire pour trouver l'expression de l'estimateur des moindres carrés mais elle est importante pour la détermination des intervalles de confiance et des p-valeurs.

6 Données du problème.

Ici, on veut expliquer les rendements de l'entreprise Total sur la période du 01/01/2015 au 31/12/2015 en fonction du :

1. cours du CAC40
2. cours du Dow Jones
3. cours d'Airfrance
4. cours de BP
5. cours de Danone
6. cours de L'Oréal
7. cours de Michelin
8. cours d'Orange
9. cours de Vinci
10. prix du baril de brut en Europe

Commentaire sur les régresseurs:

1. on peut *a priori* penser que le cours du CAC40 explique bien le cours de Total puisque les 2 sont très liés (Total étant la 1ère capitalisation boursière du CAC40)
2. on peut aussi se poser la question de l'influence des places boursières étrangères telles que l'indice Dow Jones sur le cours de Total
3. les cours d'Airfrance et Total sont sûrement anticorrélés. En effet, plus le prix du barril augmente, plus Total fait des bénéfices mais plus Airfrance doit essuyer des pertes (sa marge diminuant à cause du prix du carburant)
4. Le cours de British Petroleum doit expliquer le cours de Total (et *vice versa*)
5. *a priori* aucun lien entre Danone et Total
6. *a priori* aucun lien entre L'Oréal et Total
7. les cours de Total et Michelin doivent eu aussi être anticorrélés pour la même raison que pour Airfrance
8. *a priori* aucun lien entre Orange et Total
9. *a priori* aucun lien entre Vinci et Total
10. le cours du baril doit très bien expliquer le cours de Total (plus le premier est élevé et plus Total fait des bénéfices)

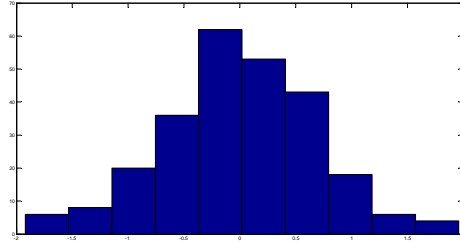
Remarque: il aurait été intéressant d'y inclure le régresseur "température" mais il est compliqué d'extraire les données d'Internet.

7 Résolution et interprétation.

D'après le cours, on sait que l'estimateur des MCO est :

$$\hat{\theta} = (\phi^T \phi)^{-1} \phi^T Y$$

Voici ci-dessous l'histogramme des résidus $e = Y - \hat{Y}$:



On constate que les résidus ont à peu près une distribution normale, ce qui est cohérent avec les conséquences des hypothèses faites au début.

7.1 Intervalles de confiance sur θ_j .

D'après le cours, on sait que

$$\frac{\hat{\theta}_j - \theta_j}{\sqrt{\hat{\sigma}^2(\phi^T \phi)^{-1}_{j,j}}} \sim Student(n - p)$$

où $\hat{\sigma}$ est l'estimateur de la variance :

$$\hat{\sigma} = \frac{\|e\|^2}{n - p}$$

d'où :

$$\mathbb{P}\left[-a \leq \frac{\hat{\theta}_j - \theta_j}{\sqrt{\hat{\sigma}^2(\phi^T \phi)^{-1}_{j,j}}} \leq a\right] = F_{Student}(a) - F_{Student}(-a)$$

ie:

$$\mathbb{P}[\theta_j \in [\hat{\theta}_j - a\sigma\sqrt{\alpha_j}, \hat{\theta}_j + a\sigma\sqrt{\alpha_j}]] = F_{Student}(a) - F_{Student}(-a), \alpha_j = (\phi^T \phi)^{-1}_{j,j}$$

On choisit a pour que le membre de droite soit égal à $1 - \alpha$ et alors $[\hat{\theta}_j - a\sigma\sqrt{\alpha_j}, \hat{\theta}_j + a\sigma\sqrt{\alpha_j}]$ est un intervalle de confiance à $100(1 - \alpha)\%$ sur θ_j . Puis en utilisant le fait que la distribution de la loi de Student est paire, on déduit qu'il faut prendre $a = -F_{Student}^{-1}(-\frac{\alpha}{2})$

Voici l'estimateur des moindres carrés et les intervalles de confiance à 95% sur θ_j :

```
theta_chapeau =
-0.0036
-0.0001
0.0771
0.0617
-0.0829
0.1083
-0.0193
0.8435
0.2562
-0.0063

intervalle =
-0.0049  -0.0022
-0.0002  0.0001
-0.1071  0.2614
0.0528   0.0706
-0.1737  0.0079
0.0738   0.1428
-0.0503  0.0118
0.6608   1.0262
0.2000   0.3124
-0.0402  0.0276
```

7.2 Sélection de variables.

Testons maintenant l'hypothèse de significativité d'un coefficient θ_j (test de Student) : on sait que

$$\frac{\hat{\theta}_j - \theta_j}{\sqrt{\hat{\sigma}^2(\phi^T \phi)^{-1}_{j,j}}} \sim Student(n-p)$$

d'où sous $H_0: \theta_j = 0$, on a : $t = \frac{\hat{\theta}_j}{\sqrt{\hat{\sigma}^2(\phi^T \phi)^{-1}_{j,j}}} \sim Student(n-p)$.

On calcule alors la p-valeur associé au test t pour chaque coefficient (qui n'est autre que $1 - F_{Student(n-p)}(t)$):

```

t =
-5.0793
-0.6891
0.8247
13.6274
-1.7993
6.1808
-1.2227
9.0936
8.9799
-0.3665

p_valeur =
1.0000
0.7543
0.2052
0
0.9634
0.0000
0.8887
0
0
0.6429

```

Or, plus la p-valeur est faible et plus on rejette H_0 l'hypothèse de nullité du coefficient. L'analyse des p-valeurs nous indique alors que :

- $\theta_4, \theta_6, \theta_8, \theta_9$ sont significatifs
- $\theta_1, \theta_2, \theta_5, \theta_7$ ne sont pas significatifs

Autrement dit, le cours de Total est bien expliqué :

- par BP (ce qui est cohérent car BP et Total sont toutes les deux des entreprises pétrolières)
- par L'Oréal (ce qui est surprenant, mais une recherche sur Internet montre qu'il y a du pétrole et du plastique dans les cosmétiques)
- le cours d'Orange et de Vinci

En revanche, le cours du CAC40, du Dow Jones, de Danone n'expliquent pas le cours de Total.

Enfin, le cours du baril de brut n'explique que moyennement le cours de Total: ceci est sûrement dû au fait que l'impact d'une montée du prix du brut n'est répercutée que plus tard dans l'année.

8 Régression Ridge.

8.1 Expression de l'estimateur Ridge.

La régression Ridge permet de corriger les 2 inconvénients suivants de l'estimateur des MCO :

- l'estimateur des MCO privilégie le biais à la variance (c'est d'ailleurs le meilleur estimateur sans biais au sens où c'est celui de variance minimale)
- il faut que $\phi^T \phi$ inversible et donc $n \geq k$, ce qui n'est pas toujours le cas

L'estimateur Ridge est alors défini comme :

$$\hat{\theta}_{Ridge} = \operatorname{argmin}(\|Y - \phi\theta\|_2^2 + \lambda\|\theta\|_2^2)$$

où λ est un paramètre de pénalisation à fixer par l'utilisateur. Nous verrons dans la suite comment le fixer de façon optimale.

D'après le cours, l'expression de $\hat{\theta}_{Ridge}$ est :

$$\hat{\theta}_{Ridge} = (\phi^T \phi + \lambda I_p)^{-1} \phi^T Y$$

8.2 Comment fixer λ : la cross-validation.

La cross-validation est une méthode pratique pour fixer le λ optimal (qui dépend de θ inconnu). On procède comme suit :

Etape 1 : on partitionne les données en K ensembles de taille égales

Etape 2 : pour chaque $k = 1, \dots, K$, on calcule $\hat{\theta}_{-k}(\lambda)$ l'estimateur Ridge sur l'ensemble des données exceptée celles du k -ième ensemble. On calcule l'erreur de prévision empirique de cet estimateur sur l'ensemble numéro k défini comme suit :

$$CVerreur_k(\lambda) = \frac{1}{\operatorname{card}(\text{ensemble}_k)} \sum_{Y_j, X_j \in \{\text{ensemble}_k\}} (Y_j - \phi_j \theta_{-k}(\lambda))^2$$

Algorithm 2 extrait de regression.m - Suite : cf annexe

```

1  for u=1:Nu
2      for k = 1:K-1
3          Y_moins_k = Y;
4          phi_moins_k = phi;
5          %Pour construire theta_ridge sur toutes les données moins celles de l'ensemble k
6          for i=1:K
7              Y_moins_k(K*(k-1)+i) = [];
8              phi_moins_k(K*(k-1)+i,:) = [];
9          end
10         %Yj, phij dans k - on ne garde que l'ensemble numéro k
11         for j=1:K
12             phi_que_k(j,:) = phi(K*(k-1)+j,:);
13             Y_que_k(j)=Y(K*(k-1)+j);
14         end
15
16         theta_ridge_moins_k = pinv((phi_moins_k'*phi_moins_k)+lambda(u)*diag(ones(p,1)))*phi_moins_k'*Y_moins_k;
17
18         V = Y_que_k'-phi_que_k*theta_ridge_moins_k;
19         CV_erreur_sur_k(k) = (1/K)*(V'*V);
20     end
21     CV_totale(u) = (1/K)*sum(CV_erreur_sur_k);
22 end

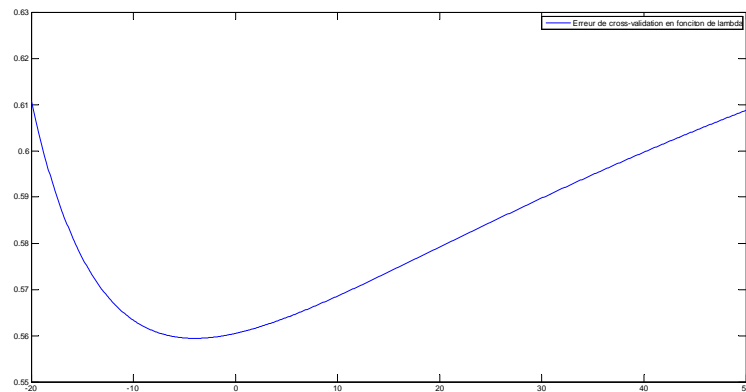
```

Etape 3 : on calcule l'erreur de cross-validation totale pour un λ donné :

$$CVerreur(\lambda) = \frac{1}{K} \sum_{k=1}^K (CVerreur_k(\lambda))$$

Enfin on choisit le λ qui minimise cette dernière quantité. Voici le script dans *regression.m* qui programme la procédure de cross-validation :

Voici le graphique représentation l'erreur de cross-validation totale en fonction de λ (le graphique en cloche est cohérent avec ce qu'on doit attendre):



On choisit donc $\lambda \approx -3.9$. L'estimateur Ridge correspondant est :

```
theta_ride_star =
-0.0037
-0.0001
0.0887
0.0612
-0.0979
0.1132
-0.0142
0.9047
0.2549
-0.0032
```

On peut alors constater que l'estimateur Ridge et l'estimateur MCO sont très proches, ce qui est cohérent (excepté pour le dernier coefficient):

```

theta_ridge_star =
-0.0037|
-0.0001      theta_chapeau =
  0.0887      -0.0036
  0.0612      -0.0001
-0.0979      0.0771
  0.1132      0.0617
-0.0142     -0.0829
  0.9047      0.1083
  0.2549     -0.0193
-0.0032      0.8435
          0.2562
          -0.0063

```

Part III

Portefeuille de Markowitz

9 Présentation du problème.

Le but d'un portefeuille de Markowitz est de trouver des combinaisons de titres (les portefeuilles) les plus intéressants statistiquement ie fournissant la meilleure performance pour un niveau de risque donné. La performance sera mesurée par la moyenne des rendements et le risque par leur variance. Autrement dit, on cherche les portefeuilles d'espérance de rendement maximal pour une variance inférieure à un niveau donné.

10 Données du problème.

Dans les 2 cas (avec ou en présence d'actif risqué), on considère un marché avec:

- une date initiale $t = 0 \stackrel{ici}{=} 02/11/2015$
- un horizon $t = 1 \stackrel{ici}{=} 22/01/2015$

Le marché est constitué de 36 actifs du CAC40 (et non 40 car certaines cotations n'étaient disponibles qu'en partie comme par exemple Alstom).

11 Cas 1 : en présence d'actif sans risque.

11.1 Mise en équation.

En plus des actifs risqués mentionnés ci-dessus, on considère un actif sans risque de taux sans risque $r = 0.03$ arbitraire (cf script *markowitz.m*)

La résolution du problème de minimisation sous contrainte :

$$\max_{a_0, a} \mathbb{E}[V_1(a_0, a)]$$

sous les contraintes :

$$\text{Var}[V_1(a_0, a)] \leq \sigma^2, V_0(a_0, a) = v$$

où :

- V : valeur du portefeuille, dépendant de
- a_0 : le nombre d'unités d'actif sans risque
- a : le nombre d'unités d'actifs risqué

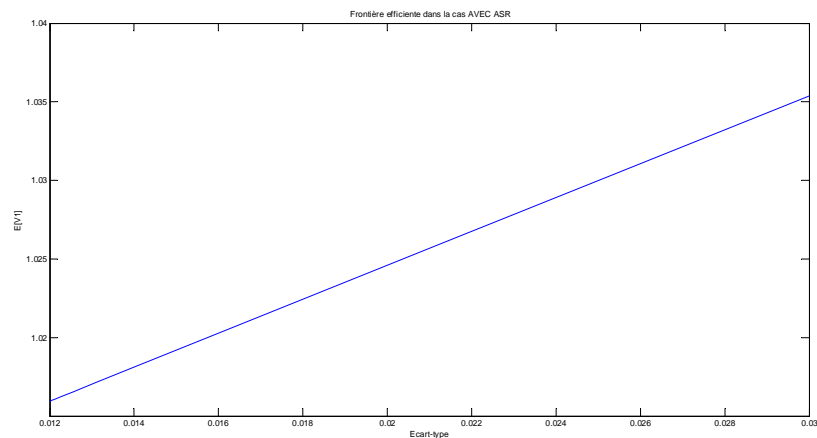
11.2 Solution.

D'après le cours, pour les portefeuilles de valeur v en $t = 0$, la frontière efficiente est constitué des portefeuilles tels que :

- $a = \frac{1}{\lambda} (\text{diag}(p_0))^{-1} \Omega^{-1} \tilde{\mu}$
- $a_0 = v - \frac{1}{\lambda} \tilde{\mu}^T \Omega^{-1} e$
- $\lambda = \frac{1}{\sigma} (\tilde{\mu}^T \Omega^{-1} \tilde{\mu})^{1/2}$

La frontière efficiente est (en paramétrisant par λ) dans le plan espérance/écart-type : $\mathbb{E}[V_1(\lambda)] = v(1 + r) + \sqrt{\text{Var}(V_1(\lambda))} (\tilde{\mu}^T \Omega^{-1} \tilde{\mu})^{1/2}$

Voici le graphique généré par *markowitz.m* pour le cas avec ASR:



```
a =
-0.0248
-0.2694
-0.0590
0.0540
-0.0645
-0.1229
0.0472
0.0118
0.0537
```

Voici un extrait du vecteur des composition :

Il faut donc détenir 0.0540 unités de l'actif 4 et vendre 0.0248 unités de l'actif 1 par exemple.

12 Cas 2 : sans ASR.

12.1 Mise en équation.

Avec les mêmes notations, le problème s'écrit :

$$\max_a \mathbb{E}[V_1(a)]$$

sous les contraintes :

$$\text{Var}[V_1(a)] \leq \sigma^2, V_0(a) = v$$

12.2 Solution.

Pour déterminer la frontière efficiente, nous allons procéder différemment. En effet, on ne dispose pas (dans le cours) de l'expression analytique de la frontière efficiente. Nous allons raisonner en % pour la composition du portefeuille ce qui sera beaucoup plus pratique. La procédure consistera alors à "plotter" un grand nombre de portefeuilles (avec des poids aléatoires en % pour chacun des actifs). La frontière efficiente se révélera d'autant mieux que le nombre de portefeuilles tirés sera grand. Puis grâce à une autre procédure, on déterminera le meilleur portefeuille pour un écart-type donné.

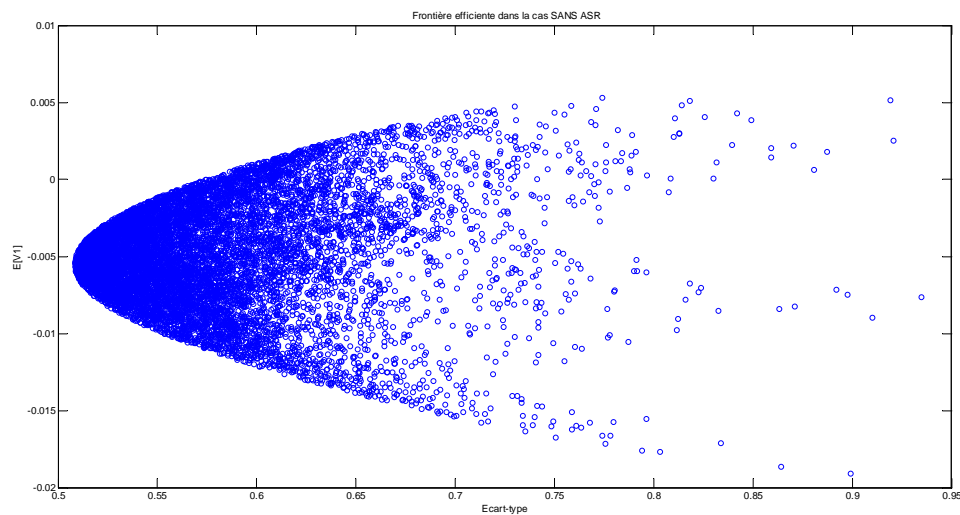
12.3 Résultat.

12.3.1 Constatation.

En considérant les rendements des 36 actifs du CAC40, on constate qu'on obtient une 1/2 parabole.

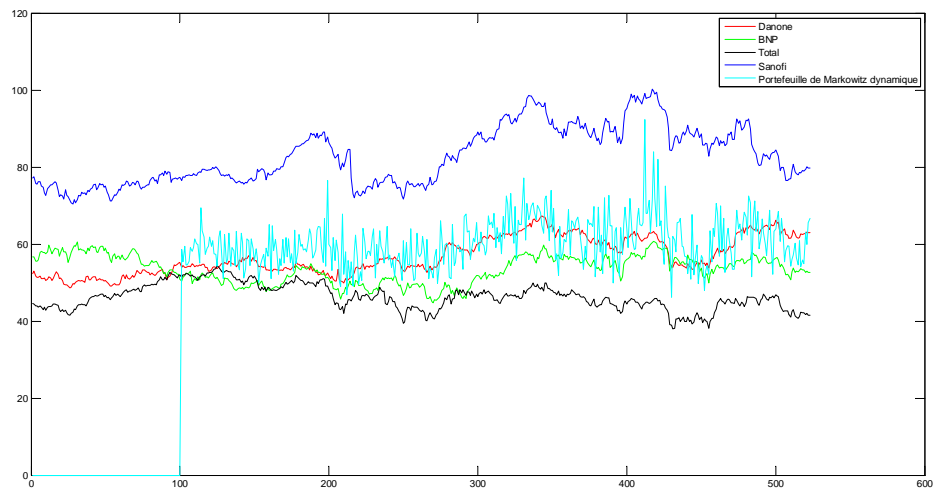
Ceci est dû au fait que la période considérée du 02/11/2015 au 22/11/2015 est trop petite, si bien que les rendements ne sont pas normaux. En effet, quand on veut tracer la frontière efficiente avec comme jeu de

données des rendements normaux, on obtient une parabole entière.



12.3.2 Exemple appliqué du trader.

Tout ceci n'est pas très appliqué. Mettons-nous alors à la place d'un trader voulant déterminer à chaque jour de trading, quel est son portefeuille de Markowitz. Considérons un jeu de données plus restreint : on va prendre seulement 4 actifs risqués (Danone, BNP, Total, Sanofi) mais sur une période plus grande du 01/01/2014 au 01/01/2016. Un des problèmes du trader va alors consister à estimer entre autre la moyenne de ses rendements: une stratégie possible est de prendre une fenêtre dynamique de 100 observations pour l'estimation des moyennes . Ces 100 observations constituent l'historique des données pour le trader. Puis, à chaque nouvelle journée de trading, cette fenêtre se décalera en intégrant le nouveau rendement (c'est une fenêtre glissante). C'est une sorte de portefeuille de Markowitz dynamique. Voici le résultat :



La procédure est assez longue (environ 4mn, mais immédiate pour le trader si il l'exécute chaque jour) car pour les 423 journées de trading, il faut à chaque fois:

- recalculer les moyennes et les covariances en intégrant le nouveau rendement
- simuler environ 1000 portefeuilles aléatoires basés sur ces rendements pour déterminer la frontière efficiente
- trouver les poids optimaux pour un niveau de risque donné

Le portefeuille performe assez bien, car il résiste bien aux périodes de baisse de 3 actifs sur 4 (période 100-250) et augmente significativement lorsque ceux-ci augmentent ensemble (période 400-425). On constate aussi que la valeur du portefeuille oscille beaucoup, ce qui est sans doute dû au fait que la moyenne des rendements estimée au fur et à mesure oscille aussi. Une manière d'améliorer l'estimation de la moyenne de ces rendements aurait été d'utiliser un filtre de Kalman.

Part IV

Annexe : code des 3 parties (ACP -
Régression MCO - Régression
Ridge - Markowitz) - cf pages
suivantes

Algorithm 3 *acp.m*

```
1 clear all;
2 close all;
3 clc;
4
5 load cac40.csv;
6
7 X = cac40;
8 [n,p] = size(cac40);
9 X = X(:,2:p);
10 [n,p] = size(X);
11
12
13 % Centrage et normalisation des points variables
14 for j=1:p
15     mu = (1/n)*sum(X(:,j)); % moyenne du point variable j
16     Y = X(:,j)-mu;
17     va = (1/n)*(Y')*Y;
18     X(:,j) = (X(:,j)-mu)/sqrt(va);
19 end
20
21 % Matrice des corrélations (matrice nxn)
22 V = (1/n)*X*X';
23
24 % Recherche des plus grandes valeurs propres et vecteurs propres
25 [vvp,vap] = eig(V); % les vecteurs propres sont normalisés
26
27
28 %%
29 % Représentation des variables sur l'axe 1
30 RepVariableAxe1 = X'*vvp(:,1); % vvp(:,n) : + grand vecteur propre de V.
31 SurAxe = zeros(1,p);
32 figure;
33 plot(RepVariableAxe1,SurAxe,'x');
34 hold on;
35 xL = xlim;
36 yL = ylim;
37 plot(xL,[0 0], 'k-');
38 legend('Représentation des variables 1 à 36 sur l'axe 1');
39 for i=1:p
40     text(RepVariableAxe1(i),SurAxe(i)+(-0.1)*mod(i,6),int2str(i))
41 end
42
43
44
45 % Projection des individus sur l'axe 1
46
47 % un individu = 1 jour de trading = 1 vecteur de taille p = 36
48 % l'axe 1 = le vecteur propre de XX' de taille (n,p)x(p,n)=(n,n) =
49 % 1 vecteur de taille n = 60
50 % La j-ième composante principale fournit les coordonnées des n individus sur le j-ième axe principal
51 % C_{j} = X*vecteur_propre_{j} de V (n,p)*(n,1)
52 figure;
53 RepIndividusAxe1 = X*(1/sqrt(vap(1,1)))*(X')*vvp(:,1);
54 plot(RepIndividusAxe1,zeros(1,n),'x');
55 legend('Projection des individus sur l'axe 1');
56 U = zeros(n,1);
57 for i=1:n
58     text(RepIndividusAxe1(i),U(i)+(-0.1)*mod(i,7),int2str(i))
59 end
60 hold on;
61 xL = xlim;
62 yL = ylim;
63 plot(xL,[0 0], 'k-');
64
65 % Représentation des variables dans le plan (axe1,axe2)
66 RepVariableAxe2 = X'*vvp(:,2);
67 figure;
68 plot(RepVariableAxe1,RepVariableAxe2,'x');
69 legend('Représentation des variables 1 à 36 dans le plan axe 1 / axe 2');
70 for i=1:p
71     text(RepVariableAxe1(i),RepVariableAxe2(i),int2str(i))
72 end
73 hold on;
74 xL = xlim;
75 yL = ylim;
76 plot(xL,[0 0], 'k-');
77 plot([0 0],yL, 'k-');
78
79 % Projection des individus dans le plan (axe1,axe2)
80 figure;
81 RepIndividusAxe2 = X*(1/sqrt(vap(2,2)))*(X')*vvp(:,2);
82 plot(RepIndividusAxe1,RepIndividusAxe2,'x');
83 legend('Projection des individus dans le plan axe 1 / axe 2');
84 for i=1:n
85     text(RepIndividusAxe1(i),RepIndividusAxe2(i)+(-0.1)*mod(i,6),int2str(i))
86 end
87 hold on;
88 xL = xlim;
89 yL = ylim;
90 plot(xL,[0 0], 'k-');
91 plot([0 0],yL, 'k-');
92
93 % Dans l'espace ...
94 RepVariableAxe3 = X'*vvp(:,3);
95 figure;
96 plot3(RepVariableAxe1,RepVariableAxe2,RepVariableAxe3,'x');
97 grid on;
98 figure;
99 RepIndividusAxe3 = X*(1/sqrt(vap(3,3)))*(X')*vvp(:,3);
100 plot3(RepIndividusAxe1,RepIndividusAxe2,RepIndividusAxe3,'x');
101 legend('Projection des individus dans l'espace axe1/axe2/axe3');
102 for i=1:n
103     text(RepIndividusAxe1(i),RepIndividusAxe2(i),RepIndividusAxe3(i),int2str(i))
104 end
105 hold on;
106 grid on;
```

Algorithm 4 *regression.m* MCO et RIDGE

```
1 clear all;
2 close all;
3 clc;
4
5 % DATES : 01/01/2015 ----- 31/12/2015 (1 an)
6 % Action Total (ou EDF) expliquée par :
7 % (100 jours de cotations par exemple = 100 obs yi)
8 % Rq : avec les cours actuels, large spectre d'observations
9 % Suivre démarche pdf régression
10
11 % ATTENTION: comparer que les jours de cotations coïncident (trous le week
12 % end)
13
14 % 1/ cours du cac40 (https://fr.finance.yahoo.com/q/hp?s=FCHI)
15 % 2/ cours du dow jones https://research.stlouisfed.org/fred2/series/DJIA/downloaddata
16 % 3/ cours d'Airfrance (qd pétrole augmente, moins de bénéf (https://fr.finance.yahoo.com/q/hp?s=AF.PA)
17 % 4/ cours d'une entreprise concurrente - BP (https://fr.finance.yahoo.com/q/hp?s=BP.L)
18 % 5/ cours d'entreprises qui n'ont (vraiment) rien à voir, par exemple
19 % 6/ cours d'entreprises qui n'ont (apparemment) rien à voir, par exemple
20 % 7/ cours de Michelin (pneus<- pétrole) (https://fr.finance.yahoo.com/q/hp?s=OR.PA)
21 % 8/ cours d'Orange https://fr.finance.yahoo.com/q/hp?s=ORA.PA&b=1&a=00&c=2015&e=31&d=11&f=2015&g=d
22 % 9/ cours Vinci https://fr.finance.yahoo.com/q/hp?s=DG.PA&b=1&a=00&c=2015&e=31&d=11&f=2015&g=d
23 % 10/ prix du baril de brut Europe (taper : https://research.stlouisfed.org/fred2/series/DCOILBRENTU/downloaddata)
24
25
26 % de 01/01/2010 à 01/01/2015, chaque jour
27
28 % Régression linéaire simple
29 % Variable observée (cours d'ouverture de total (2eme colonne))
30 load total.csv
31
32 % Variables explicatives
33 load cac40.csv; load dowjones.csv; load af.csv; load bp.csv; load danone.csv; load loreal.csv; load michelin.csv; load orange.csv; load vinci.csv; load brut.csv
34
35 n=size(total(:,2));
36 n=n(1);
37 p = 10; % nombre de régresseurs
38
39 % Estimateurs MC
40 phi = [cac40(:,2), dowjones(:,2), af(:,2), bp(:,2), danone(:,2), loreal(:,2), michelin(:,2), orange(:,2), vinci(:,2), brut(:,2)];
41 theta_chapeau = pinv(phi'*phi)*phi'*total(:,2);
42 Ychapeau = phi*theta_chapeau;
43
44 % Histogramme des résidus
45 Y = total(:,2);
46 e = Y-Ychapeau;
47 hist(e);
48
49 % Intervalles de confiance à 95% sur theta_j
50 e = Y-Ychapeau;
51 MSE = e'*e/(n-p); %non n-p+1
52 sigma_hat = sqrt(MSE);
53 alpha = 0.05;
54 quantile = -qt(alpha/2,n-p); %qt: quantile de la loi de Student à n-(p+1) degrés de liberté
55 intervalle = zeros(p,2);
56 M = pinv(phi'*phi);
57 for j=1:p
58     a = quantile*sigma_hat*sqrt(M(j,j));
59     intervalle(j,:) = theta_chapeau(j)+[-a,a];
60 end
61 intervalle
62
63 % Sélection de variables
64 % p-valeurs
65 std = sqrt(diag(MSE*inv(phi'*phi)));
66 t = theta_chapeau./std % fournit un test de student pour chaque coefficient
67 p_valeur = 1-tcdf(t,n-p)
68 % Plus la p-valeur est faible, plus on a confiance dans H1 ie plus on
69 % rejette H0: l'hypothèse de nullité du coefficient
70
71 % Régression Ridge
72 Y = total(:,2);
73 phi = [cac40(:,2), dowjones(:,2), af(:,2), bp(:,2), danone(:,2), loreal(:,2), michelin(:,2), orange(:,2), vinci(:,2), brut(:,2)];
74 lam = 2;
75 theta_ridge = pinv((phi'*phi)+lam*diag(ones(p,1)))*phi'*Y;
76 % Choix de lambda - Cross validation: ici n = 256: on partitionne en K = 16
77 % ensembles de 16 éléments chacun
78 K = 16;
79 lambda = [-20:0.1:50];
80 Nu = size(lambda);
81 Nu = Nu(2);
82
83 for u=1:Nu
84     for k = 1:K-1
85         Y_moins_k = Y;
86         phi_moins_k = phi;
87         %Pour construire theta_ridge sur toutes les données moins celles de l'ensemble k
88         for i=1:K
89             Y_moins_k(K*(k-1)+i) = [];
90             phi_moins_k(K*(k-1)+i,:) = [];
91         end
92         %Yj, phi_j dans k - on ne garde que l'ensemble numéro k
93         for j=1:K
94             phi_que_k(j,:) = phi(K*(k-1)+j,:);
95             Y_que_k(j)=Y(K*(k-1)+j);
96         end
97
98         theta_ridge_moins_k = pinv((phi_moins_k'*phi_moins_k)+lambda(u)*diag(ones(p,1)))*phi_moins_k'*Y_moins_k;
99
100         V = Y_que_k'-phi_que_k*theta_ridge_moins_k;
101         CV_erreur_sur_k(k) = (1/K)*(V'*V);
102     end
103     CV_totale(u) = (1/K)*sum(CV_erreur_sur_k);
104 end
105 figure;
106 plot(lambda,CV_totale)
107 legend('Erreur de cross-validation en fonction de lambda');
108 [CV_totale_min,u_cvmin] = min(CV_totale)
109 lambda_star = lambda(1) + 0.1*u_cvmin;
110 theta_ridge_star = pinv((phi'*phi)+lambda_star*diag(ones(p,1)))*phi'*Y
```

Algorithm 5 *markowitz.m* pour la frontière efficiente avec ASR

```

1  clear all;
2  close all;
3  clc;
4
5  % Source : Yahoo finance
6  % 36 actifs (et non 40 car certaines cotations n'étant disponibles qu'en partie comme par exemple Alstom
7  % du 02/11/2015 au 22/01/2016
8
9  load cac40.csv;
10
11 % n = nombre d'obs / p = nb d'actifs + 1
12 c = cac40;
13 [n0,p0] = size(c);
14 prix = cac40(:,2:p0);
15 [n,p] = size(prix);
16
17 % Vecteur des "rendements" entre t0 = 21/01/2016 et t = 22/01/2016
18 yi = prix(n,:)./prix(n-1,:);
19
20 % Vecteur des moyennes de rendements
21 y = (prix(2:n,:)./prix(1:n-1,:)); % Attention hypothèse  $\mu_i > 1+r$ 
22 mu_hat = mean(y) ;
23
24 % Vecteur des moyennes des excédents de rendement
25 r = 0.003;
26 mu_tilde = (mu_hat -(1+r)*ones(1,p))';
27
28
29 % Matrice de corrélation empirique
30
31 % Vecteur colonne contenant les rendements de chaque actif i = 1,...,p
32 yi = yi';
33 mu_hat = mu_hat';
34
35 % Estimateur empirique de la matrice de corrélation
36 %omega = (1/n)*(yi-mu_hat)*(yi-mu_hat)'; <- MAUVAIS ESTIMATEUR
37 omega = cov(y); % <- MEILLEUR ESTIMATEUR (MATLAB)
38
39
40 % Matrice dont la diagonale est p_i,0 (actif i, date t = 0 = 21/01/2016 ie n-1-ième ligne du tableau
41 diag_po = diag(prix(n-1,:));
42
43 % CAS 1 : AVEC ASR
44 % Portefeuille optimal pour sigma et v donné arbitrairement...
45
46 sigma = 0.1; % Niveau de risque
47 v = 1; % Valeur initiale du portefeuille
48
49 la = sqrt(mu_tilde'*pinv(omega)*mu_tilde)/sigma
50 ao = v -(mu_tilde'*pinv(omega)*ones(p,1)/la)
51 a = pinv(diag_po)*pinv(omega)*mu_tilde/la %
52
53 % Frontière optimale
54
55 %%for k = 1:10
56 %sig = sqrt(mu_tilde'*pinv(omega)*mu_tilde)/lambda
57 %%% EVlambdak = v*(1+r)+(k-1)*sqrt(mu_tilde'*pinv(omega)*mu_tilde);
58 %%%end
59 ma_fonction = @(ecart_type_portefeuille) v*(1+r)+ ecart_type_portefeuille*sqrt(mu_tilde'*pinv(omega)*mu_tilde);
60 fplot(ma_fonction,[0.012,0.03]);
61 %plot([0:9],EVlambdak);
62 title('Frontière efficiente dans la cas AVEC ASR');
63 xlabel('Ecart-type');
64 ylabel('E[V1]');
65
66 % CAS 1 : SANS ASR
67 % Frontière optimale
68
69
70 hold on;
71 figure
72 M = 10000;
73
74
75 % Valeur poretfeuille = Pt = SUM(ponderations_i x Valeurs_actif_i)
76 % On raisonne en pourcentage
77
78 %rend = [0.05 0.08]';
79 %m_o = [0.03 0.3 ; 0.3 0.05];
80 %omega(1:2,1:2)
81 %omega = [0.04 0.001 0.002 ; 0.001 0.5 0.06; 0.002 0.06 0.001];
82 %omega = [0.04 0.01 ; 0.01 0.5];
83
84 for i = 1:M
85 %ponderation = randn(p,1); ATTENTION : ce sont les rendements qui sont
86 %normaux, pas les poids. Il y a symétrie dans la parabole si les
87 %rendements sont normaux. Donc prendre plus de données, pas que 60
88 %journées mais sur 1 ans et avec quelques actifs seulement
89 ponderation = rand(p,1);
90 w = ponderation./sum(ponderation); %poids en %
91 esperance = (w')*mu_hat;
92 %omega = single(omega);
93 variance = (w')*omega*w; % normal si 24
94 % Pour éviter les points extrêmes (pour le plot)
95 if (sqrt(variance)<0.5)
96 plot(sqrt(variance),esperance,'o');
97 end
98 hold on;
99 end
100 title('Frontière efficiente dans la cas SANS ASR');
101 xlabel('Ecart-type');
102 ylabel('E[V1]');

```

Algorithm 6 *markowitzbis.m* frontière efficiente dans le cas sans ASR et avec des rendements normaux

```

1  clear all;
2  close all;
3  clc;
4
5  p = 4 ; %nombres d'actifs
6  n = 1000 ; %nombre de jours de trading
7  M = 10000; % nombre de portefeuilles aléatoires
8  %ATTENTION : il y a symétrie dans la parabole si les rendements sont normaux. Donc prendre plus de données ...
9
10 rendements = randn(n,p);
11 mu_hat = (mean(rendements))';
12 omega = cov(rendements);
13
14 risque = 0.6;
15 epsilon = 0.01;
16 res = -15;
17 % Tracer la frontière efficiente
18 % Et pour un niveau de risque donné, donne l'espérance de rendement
19 % correspondante ainsi que la pondération associée
20 for i = 1:M
21     aux = 1;
22     ponderation = rand(p,1);
23     w = ponderation./sum(ponderation); %poids en %
24     esperance = (w')*mu_hat;
25     variance = (w')*omega*w; % il faut omega dans Sn++
26     if (sqrt(variance)<1) % pour éviter les points extrêmes (pour le plot)
27         plot(sqrt(variance),esperance,'o');
28     end
29     if (abs(sqrt(variance)-risque)<epsilon) && (esperance>res)
30         res = esperance;
31         w_star = w;
32     end
33     hold on;
34 end
35 res
36 title('Frontière efficiente dans la cas SANS ASR');
37 xlabel('Ecart-type');
38 ylabel('E[V1]');

```

Algorithm 7 *markody.m* fonction appelée à chaque nouveau jour de trading

```

1  function [w_star,res,bool] = marko_dy(p,M,rendements,risque,epsilon)
2  mu_hat = (mean(rendements))';
3  omega = cov(rendements);
4
5  ponderation = rand(p,1);
6  w_star = ponderation./sum(ponderation);
7  w_star_init = w_star;
8  res = -15;
9
10 % Tracer la frontière efficiente
11 % Et pour un niveau de risque donné, donne l'espérance de rendement
12 % correspondante ainsi que la pondération associée
13 %w_star = zeros(p,1); % au pire tt faire ds même fenêtre
14 for i = 1:M
15     ponderation = rand(p,1);
16     w = ponderation./sum(ponderation); %poids en %
17     esperance = (w')*mu_hat;
18     variance = (w')*omega*w; % il faut omega dans Sn++
19     if (sqrt(variance)<1) % pour éviter les points extrêmes (pour le plot)
20         %plot(sqrt(variance),esperance,'o'); <- pas besoin de ploter
21         %ici
22     end
23     if (abs(sqrt(variance)-risque)<epsilon) && (esperance>res) %Pour trouver w_star pour un niveau de risque c
24         res = esperance;
25         w_star = w;
26     end
27     hold on;
28 end
29 %title('Frontière efficiente dans la cas SANS ASR');
30 %xlabel('Ecart-type');
31 %ylabel('E[V1]');
32 bool = (w_star_init==w_star); % Pour vérifier que ça a bien optimisé...

```

Algorithm 8 *markodynamique.m* script plotant la valeur du portefeuille de Markowitz au cours du temps et appelant la fonction *markody.m*

```

1  clear all;
2  close all;
3  clc;
4
5  % 01/01/2014 — 01/01/2016
6  load actifs.csv;
7
8  c = actifs;
9  [n0,p0] = size(c); % n0 = nombre d'obs / p0 = nb d'actifs + 1
10 prix = actifs(:,2:p0);
11 [n,p] = size(prix); % n = nombre d'observations TOTAL, p = nombre d'actifs
12
13 figure;
14 plot(1:n, prix(1:end,1), 'r');
15 hold on;
16 plot(prix(1:end,2), 'g');
17 hold on;
18 plot(prix(1:end,3), 'k');
19 hold on;
20 plot(prix(1:end,4), 'b');
21 hold on;
22
23 % rendements = prix(n,:)./prix(n-1,:); % vecteur des "rendements" entre
24 % t0 = 21/01/2016 et t = 22/01/2016
25 % le fait que la frontière ne soit pas belle est dû au fait que les
26 % rendements ne sont pas normaux. En effet, cf figure où la frontière est
27 % "belle" dans le cas de rendements normaux...
28
29 % On se place dans le rôle du trader qui voit arriver les cotations au fur
30 % et à mesure.
31 % On prend une fenêtre dynamique de 100 observations pour l'estimation de mu
32 % Ces 100 observations constituent l'historique des données pour le trader
33 %for i = 101:n
34
35 M = 100;
36 risque = 0.012;
37 epsilon = 0.00001;
38 fen = 100;
39 for i = 101:523
40     rendements = (prix(2+(i-101):fen+(i-101),:)./prix(1+(i-101):fen+(i-101)-1,:));
41     [w_star,res,bool] = marko_dy(p,M,rendements,risque,epsilon);
42     cotations = (prix(i,:))';
43     portefeuille(i) = (w_star)*cotations;
44 end
45 %size([101:300])
46 %size(portefeuille)
47 plot([1:523],portefeuille,'-c');
48 legend('Danone','BNP','Total','Sanofi','Portefeuille de Markowitz dynamique');

```
