



# **Apprentissage Statistique Automatique I**

## Régression linéaire

---

Andrés F. López-Lopera  
Université Polytechnique Hauts-de-France (UPHF)

## 1. Régression linéaire

Régression linéaire simple

Régression linéaire multiple

## 2. Décomposition en valeurs singulières

Définition

SVD et moindres carrés

Analyse du biais, de la variance et du risque par la SVD

Stabilité numérique

## 3. Réduction de la dimension

Régression sur les composantes principales

Régression des moindres carrés partiels

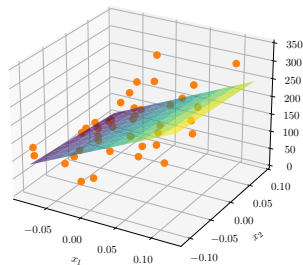
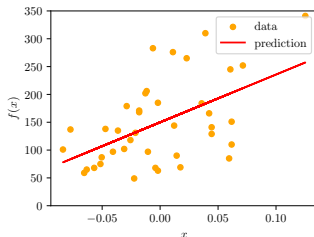
## Régression linéaire

---

- Soit

$$f : \begin{cases} \mathbb{R}^d \rightarrow \mathbb{R} \\ x_1, \dots, x_d \mapsto f(x_1, \dots, x_d) \end{cases}$$

- La **régression linéaire** est une approche simple pour l'apprentissage supervisé, supposant que  $y$  dépend linéairement de  $x_1, \dots, x_d$



- Bien que simple, elle est utile à la fois conceptuellement et en pratique

- Dans un cadre d'apprentissage, on cherche à créer une règle de régression dans la classe

$$\mathcal{F}_L := \left\{ f(\mathbf{x}) = \beta_0 + \sum_{j=1}^d \beta_j x_j, \beta = (\beta_0, \dots, \beta_d) \in \mathbb{R}^{d+1} \right\}, \quad (1)$$

avec  $\mathbf{x} = (x_1, \dots, x_d)^\top \in \mathbb{R}^d$

- On supposera qu'il existe  $\beta$  (inconnu) t.q.

$$y(\mathbf{x}) = \beta_0 + \sum_{j=1}^d \beta_j x_j + \varepsilon, \quad (2)$$

où  $\varepsilon$  est un bruit additif centré, *i.e.*,  $\mathbb{E}(\varepsilon) = 0$ . Par exemple,  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

- On pose un modèle de la forme

$$y(x) = \beta_0 + \beta_1 x + \varepsilon \quad (3)$$

- Le modèle précédent correspond bien à l'équation d'une droite avec
  - $\beta_0$  : ordonnée à l'origine (*intercept*)
  - $\beta_1$  : pente (*slope*)
- $\beta = (\beta_0, \beta_1)$  sont les **coefficients** du **modèle linéaire**
- Étant estimé  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)$ , on peut ensuite prédire  $y$  pour une valeur de  $x$

$$\hat{y}(x) = \hat{\beta}_0 + \hat{\beta}_1 x \quad (4)$$

# Méthode de moindres carrés

- Supposons une base de données  $(x, y) = (x_i, y_i)_{1 \leq i \leq n}$  avec  $x_i \in \mathbb{R}$  et  $y_i := f(x_i) \in \mathbb{R}$
- Soit  $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  la prédiction de  $y$  sur la  $i$ -ème valeur de  $X$
- Pour estimer  $(\beta_0, \beta_1)$ , on cherche à minimiser la somme des carrés résiduels (SSE, en anglais)

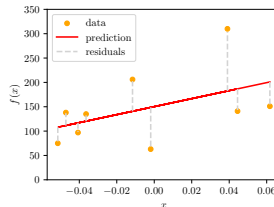
$$\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin}_{\beta \in \mathbb{R}^2} \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (5)$$

ce qui donne **[exercice]**

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\operatorname{cov}(x, y)}{\operatorname{var}(x)}$$

$$\text{avec } \bar{z} = \frac{1}{n} \sum_{i=1}^n z_i$$



**Solution.** Conditions nécessaires d'optimalité (CNO) du premier ordre :

$$\cdot \frac{\partial \text{SSE}}{\partial \beta_0} = 0 :$$

$$\frac{\partial \text{SSE}}{\partial \beta_0} = \frac{1}{2} \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = - \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0,$$

d'où on obtient

$$\hat{\beta}_0 = \frac{1}{n} \left[ \sum_{i=1}^n (y_i - \beta_1 x_i) \right] = \frac{1}{n} \sum_{i=1}^n y_i - \beta_1 \frac{1}{n} \sum_{i=1}^n x_i = \bar{y} - \beta_1 \bar{x} \quad (6)$$

$$\cdot \frac{\partial \text{SSE}}{\partial \beta_1} = 0. \text{ En sachant (6),}$$

$$\frac{\partial \text{SSE}}{\partial \beta_1} = \frac{1}{2} \frac{\partial}{\partial \beta_1} \sum_{i=1}^n (y_i - \hat{\beta}_0 - \beta_1 x_i)^2 = \sum_{i=1}^n (y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i)(\bar{x} - x_i) = 0,$$

d'où on obtient

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



## Erreur Quadratique Moyenne (*Mean Squared Error, MSE*) :

$$MSE = \frac{SSE}{n}, \quad \text{avec } SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (7)$$

## Erreur Quadratique Moyenne Standardisée (*Standardized MSE, SMSE*) :

$$SMSE = \frac{MSE}{SST}, \quad \text{avec } SST = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2, \quad (8)$$

où SST est la somme des carrés totale

## Fraction de Variance Expliquée :

$$R^2 = \frac{SST - MSE}{SST} = 1 - SMSE \quad (9)$$

Le  $R^2$  est également connu comme le coefficient de corrélation entre Y et X

**Rem :** Le  $R^2$  devient  $Q^2$  s'il est calculé que sur une base de test

## F-statistique :

$$F_{stat} = \left( \frac{n - d - 1}{d} \right) \frac{SST - SSE}{SSE} = \left( \frac{n - d - 1}{d} \right) \frac{R^2}{1 - R^2} \quad (10)$$

## sklearn.metrics

- SSE : mean\_squared\_error
- $R^2$  : r2\_score

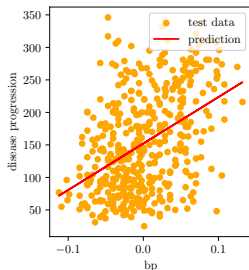
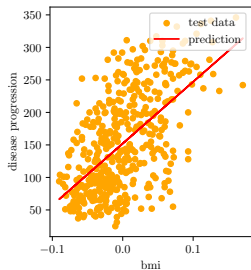
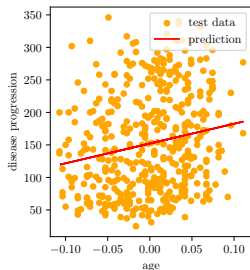
## sklearn.feature\_selection

- $F_{stat}$ ,  $p_{val}$  : f\_regression

Plus de détails :

- [https://scikit-learn.org/1.5/modules/model\\_evaluation.html](https://scikit-learn.org/1.5/modules/model_evaluation.html)
- [https://scikit-learn.org/1.5/modules/feature\\_selection.html](https://scikit-learn.org/1.5/modules/feature_selection.html)

Caractéristiques :  $n = 442$ ,  $d = 3$



$x_i$	$\beta_0$	$\beta_1$	MSE ( $\downarrow$ )	SMSE ( $\downarrow$ )	$R^2$ ( $\rightarrow 1$ )	$F_{stat}$ ( $\uparrow$ )	$P_{value}$
âge	152.13	304.18	5720.55	0.93	0.04	1.58	$7.06 \times 10^{-5}$
bmi	152.13	949.44	3890.46	0.66	0.34	22.59	$3.47 \times 10^{-42}$
bp	152.13	714.74	4774.10	0.81	0.19	10.43	$1.65 \times 10^{-22}$

- bmi : indice de masse corporelle (*body mass index*)
- bp : tension artérielle (*blood pressure*)

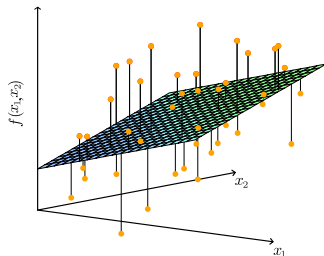
- On pose un modèle de la forme

$$y(x) = \beta_0 + \sum_{j=1}^d \beta_j x_j + \varepsilon, \quad (11)$$

avec  $x = (x_1, \dots, x_d) \in \mathbb{R}^d$

- A partir de  $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_d)$ , on peut ensuite prédire  $y$  avec

$$\hat{y}(x) = \hat{\beta}_0 + \sum_{j=1}^d \hat{\beta}_j x_j \quad (12)$$



- Comme pour  $d = 1$ , on estime  $\beta$  en minimisant la  $SSE = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$
- L'estimation dépend d'une inversion de matrice produit qui est possible de calculer en utilisant un logiciel de statistique !
- De même,  $\text{var}(\hat{\beta}_j)$ , t-test de nullité, test de Fisher, ...

# Régression linéaire multiple

- Pour une démonstration rapide, considérons une base de données  $(X, y) = (x_i, y_i)_{1 \leq i \leq n}$  avec  $x_i = (x_{i,1}, \dots, x_{i,d}) \in \mathbb{R}^d$  et  $y_i \in \mathbb{R}$
- De façon matricielle, on a  $y = X\beta$  avec

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_d \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,d} \end{bmatrix}.$$

- Alors, par minimisation du critère des moindres carrés :

$$\sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{i,j} \right)^2 = \|y - X\beta\|^2 = y^\top y - 2\beta^\top X^\top y + \beta^\top X^\top X \beta$$

- Par dérivation matricielle (sous l'hypothèse  $X^\top X$  est une matrice de plein rang), on obtient :

$$\hat{\beta} = (X^\top X)^{-1} X^\top y,$$

dont la solution correspond bien à un minimum car la matrice hessienne  $2X^\top X$  est semi définie-positive

**Caractéristiques :**  $n = 442$ ,  $d = 3$

Variable	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	MSE	SMSE	$Q^2$
(âge,bmi,bp)	152.13	25.99	788.78	394.13	<b>3580.33</b>	<b>0.60</b>	<b>0.40</b>
âge	152.13	304.18	–	–	5720.55	0.93	0.04
bmi	152.13	–	949.44	–	3890.46	0.66	0.34
bp	152.13	–	–	714.74	4774.10	0.81	0.19

- bmi : indice de masse corporelle (*body mass index*)
- bp : tension artérielle (*blood pressure*)

· Que pourrait-on conclure du tableau ?

**Caractéristiques :**  $n = 442, d = 3$

Variable	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	MSE	SMSE	$Q^2$
(âge,bmi,bp)	152.13	25.99	788.78	394.13	<b>3580.33</b>	<b>0.60</b>	<b>0.40</b>
âge	152.13	304.18	–	–	5720.55	0.93	0.04
bmi	152.13	–	949.44	–	3890.46	0.66	0.34
bp	152.13	–	–	714.74	4774.10	0.81	0.19

- bmi : indice de masse corporelle (*body mass index*)
- bp : tension artérielle (*blood pressure*)

• Que pourrait-on conclure du tableau ?

- Interprétation de  $\beta_j$  comme l'effet moyen sur  $y$  d'un accroissement de  $x_j$  d'une unité (lorsque tous les autres prédicteurs sont fixés)
- Impossible de faire des affirmation en terme de **causalité**. Comment doit-il s'interpréter  $\beta_1 = 25.99$  ?

1. Y a t-il au moins un des  $x_1, \dots, x_d$  utile pour prédire  $y$  ?
2. Sont-ils vraiment tous utiles ?
3. Comment le modèle s'ajuste aux données ?
4. Avec une nouvelle valeur de  $x_*$ , quelle réponse doit-on prédire ?  
Précision de la prédiction ?



1. Y a t-il au moins un des  $x_1, \dots, x_d$  utile pour prédire  $y$  ?
  3. Comment le modèle s'ajuste aux données ?
- On peut utiliser la  $F$ -statistique (et/ou les autres critères d'évaluation de la précision)

Variable	MSE ( $\downarrow$ )	SMSE ( $\downarrow$ )	$Q^2$ ( $\rightarrow 1$ )	$F_{stat}$ ( $\uparrow$ )
(âge,bmi,bp)	<b>3580.33</b>	<b>0.60</b>	<b>0.40</b>	<b>28.28</b>
âge	5720.55	0.93	0.04	1.58
bmi	3890.46	0.66	0.34	22.59
bp	4774.10	0.81	0.19	10.43

Cas de test - diabète (sklearn)

## 2. Sont-ils vraiment tous utiles ?

### Choix de co-variables

- I. Approche complète : Comparer les modèles linéaires avec tous les sous-ensembles possibles de co-variables (Souvent  $2^p$  trop grand)
- II. Approche séquentielle que ne parcourt que certains sous-ensembles.

---

#### Algorithm Sélection progressive

---

**Données d'entrée :**  $(X, y)$ ,  $\text{tol } \epsilon > 0$

- 1:  $D^{(0)} = \{1, \dots, d\}, X_* = ()$
  - 2: **for**  $i = 1, \dots, d$  **do**
  - 3:     **for**  $j \in D_{i-1}$  **do**
  - 4:          $e_j^2 = \text{SSE}(\mathcal{F}_L(X_* \cup_c X_{:,j}), y)$
  - 5:      $j_* = \underset{j \in D_{i-1}}{\text{argmin}} e_j^2$
  - 6:      $D^{(i)} = D^{(i-1)} \setminus j_*$
  - 7:      $X_* = X_* \cup_c X_{:,j_*}$
  - 8:     **if** critère  $\leq \epsilon$  **then**
  - 9:         End
- 

---

#### Algorithm Sélection rétrograde

---

**Données d'entrée :**  $(X, y)$ ,  $\text{tol } \epsilon > 0$

- 1:  $D^{(0)} = \{1, \dots, d\}$
  - 2: **for**  $i = 1, \dots, d$  **do**
  - 3:     **for**  $j \in D_{i-1}$  **do**
  - 4:          $D_{-j} = D^{(i-1)} \setminus j$
  - 5:          $p_j = p\text{-val}(\mathcal{F}_L(X_{:,D_{-j}}), y)$
  - 6:      $j_* = \underset{j \in D_{i-1}}{\text{argmax}} p_j$
  - 7:      $D^{(i)} = D^{(i-1)} \setminus j_*$
  - 8:      $X = X_{:,D^{(i)}}$
  - 9:     **if** critère  $\leq \epsilon$  **then**
  - 10:         End
-

4. Avec une nouvelle valeur de  $x_*$ , quelle réponse doit-on prédire ?  
Précision de la prédiction ?

- Prédiction :

$$\hat{y}(x) = \hat{\beta}_0 + \sum_{j=1}^d \hat{\beta}_j x_j,$$

pour tout  $x \in \mathbb{R}^d$  !

- Précision de la prédiction : MSE, SMSE,  $Q^2$ , ...

**Variables qualitatives.** On crée des nouvelles variables binaires (autant que nécessaire pour représenter les modalités) :

$$x_j = \begin{cases} 1 & \text{si la condition est satisfaite,} \\ 0 & \text{sinon} \end{cases}$$

**Interactions de variables.** On crée des nouvelles variables données par des produits de variables d'intérêt :

$$y(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 \underbrace{(x_1 x_2)}_{x_3}$$

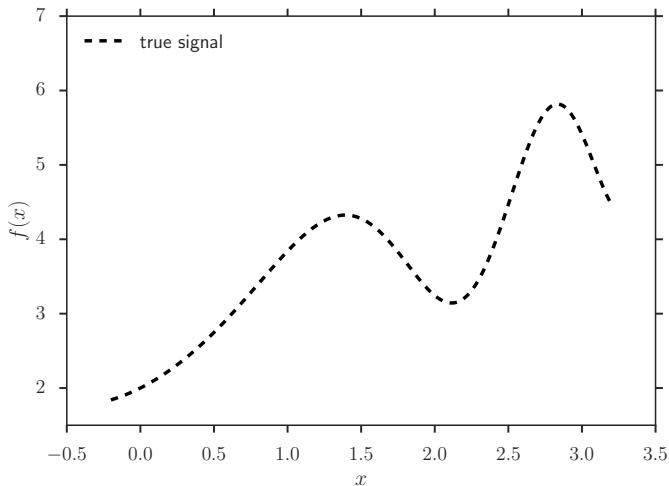
**Effets non linéaires.** On crée des nouvelles variables données par la puissance de la variable d'intérêt :

$$y(x) = \beta_0 + \sum_{j=1}^p \beta_j x^j + \varepsilon \quad (\text{régression polynomiale d'ordre } p)$$

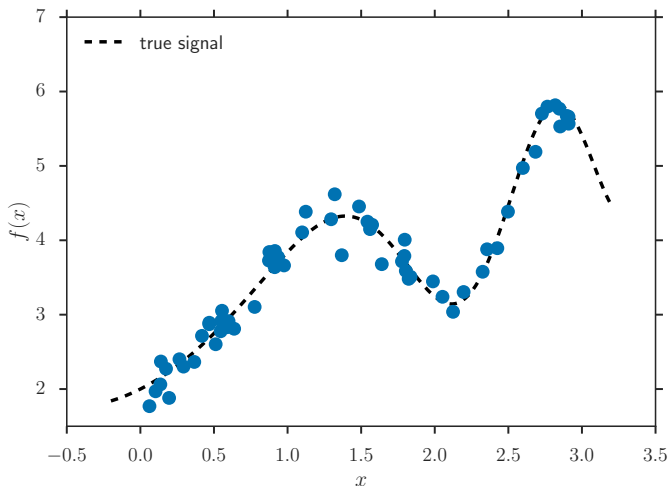
En général, on peut envisager une projection sur un ensemble de fonctions de base  $\phi_1, \dots, \phi_p : \mathbb{R}^d \rightarrow \mathbb{R}$

$$y(x) = \beta_0 + \sum_{j=1}^p \beta_j \phi_j(x) + \varepsilon$$

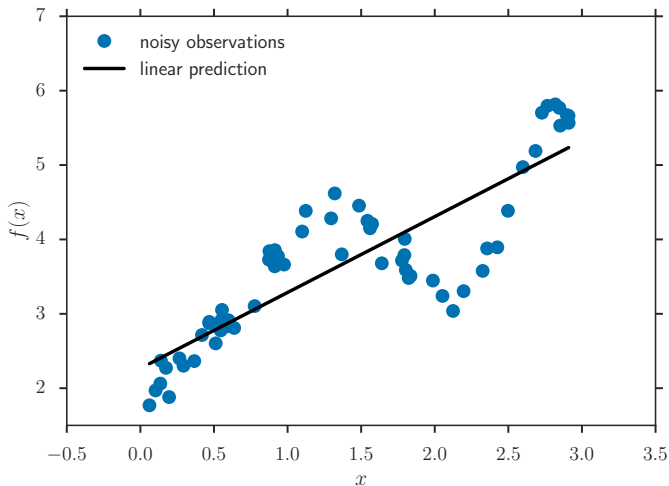
Fonction cible :  $f(x_i)$  for  $i = 1, \dots, n$



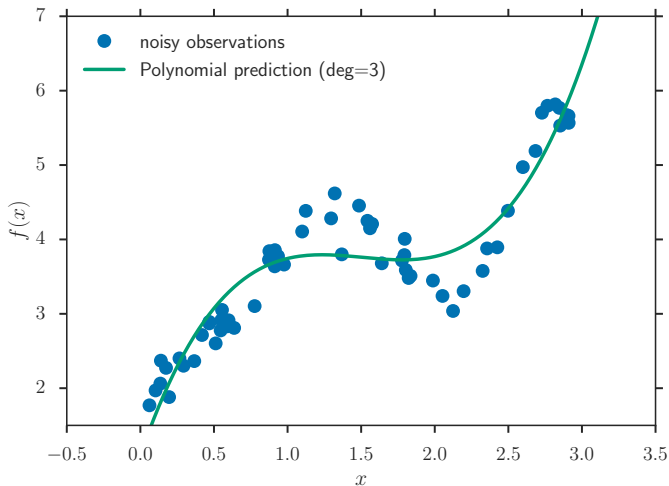
Observations bruitées :  $y_i = f(x_i) + \varepsilon_i$  pour tout  $i = 1, \dots, n$



## Régression linéaire

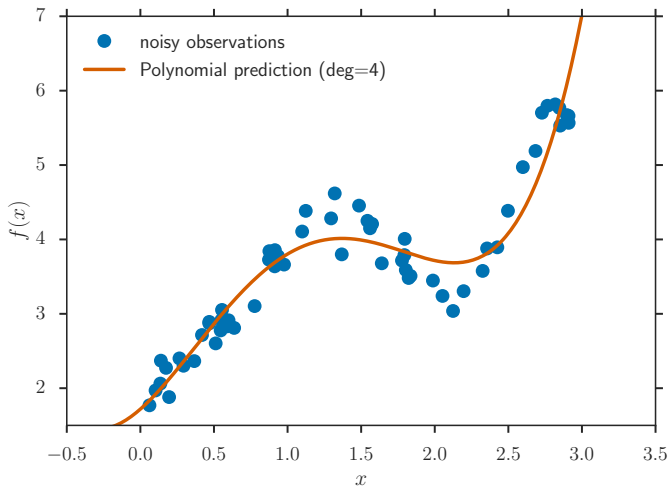


## Régression polynomiale

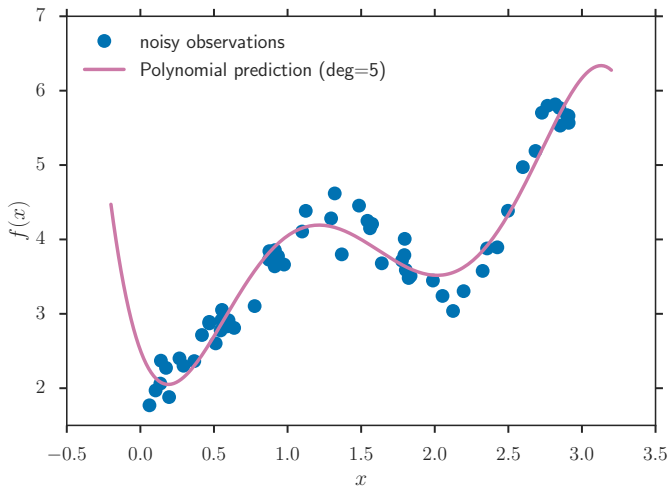




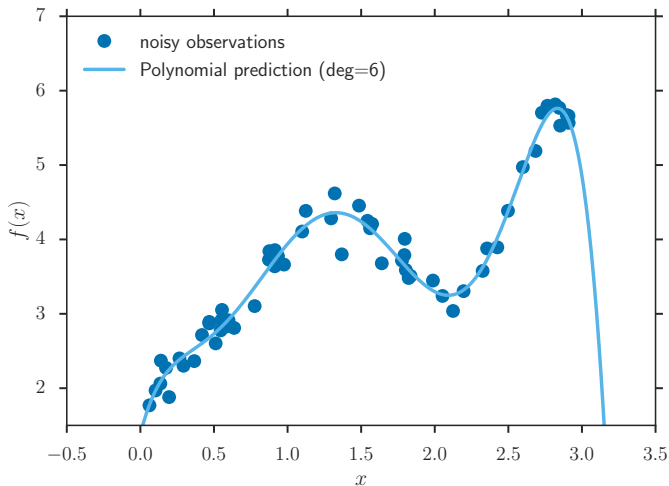
## Régression polynomiale



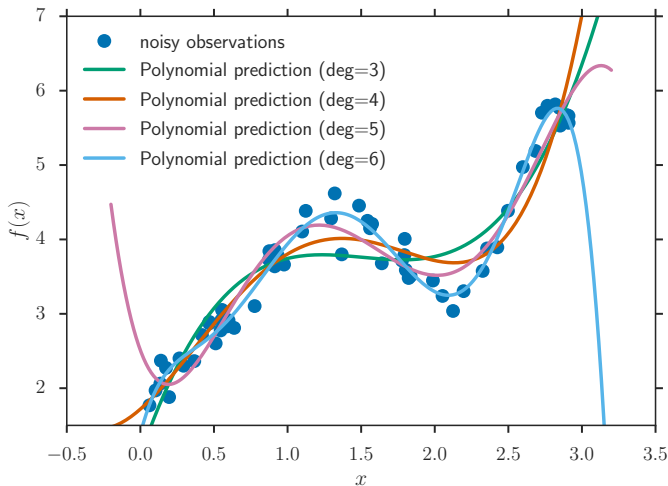
## Régression polynomiale



## Régression polynomiale



## Régression polynomiale



- **Ajustement régularisé** : régression Ridge, Lasso, Elastic Net, ...
- **Problèmes de classification** : régression logistique, support vector machine
- **Non-linéarité** : lissage à noyau, splines, modèles additifs généralisés, méthode de plus proches voisins
- **Interactions** : méthode basée sur des arbres, bagging, forêts aléatoires (random forests), et boosting (capturent aussi les non-linéarités)

## Décomposition en valeurs singulières

---

## Theorem (Théorème spectral [Golub and Van Loan, 2013])

Une matrice symétrique  $S \in \mathbb{R}^{n \times n}$  est diagonalisable en base orthonormée, i.e., il existe  $\lambda_1 \geq \dots \geq \lambda_n$  et une matrice orthogonale  $U \in \mathbb{R}^{n \times n}$  t.q. :

$$S = U \Sigma U^T \Leftrightarrow SU = U \Sigma,$$

avec  $\Sigma = \text{diag}(\lambda_1, \dots, \lambda_n)$  et  $\lambda_1 \geq \dots \geq \lambda_n > 0$

· Si l'on écrit  $U = [u_1, \dots, u_n]$  cela signifie que :

$$S = \sum_{i=1}^n \lambda_i u_i u_i^T, \quad \text{avec} \quad S u_i = \lambda_i u_i \text{ pour tout } i \in \llbracket 1, n \rrbracket$$

Rappel : une matrice orthogonale  $U \in \mathbb{R}^{n \times n}$  doit vérifier :

$$U^T U = U U^T = I \Leftrightarrow u_i^T u_j = \langle u_i, u_j \rangle = \delta_{i,j}$$

· Les  $\lambda_i \in \mathbb{R}^+$  sont les **valeurs propres** de  $S$  et les  $u_i \in \mathbb{R}^n$  sont les **vecteurs propres** associés

## Theorem (Décomposition SVD [Golub and Van Loan, 2013])

Pour toute matrice  $X \in \mathbb{R}^{n \times p}$ , il existe deux matrices orthogonales  $U = [u_1, \dots, u_n] \in \mathbb{R}^{n \times n}$  et  $V = [v_1, \dots, v_p] \in \mathbb{R}^{p \times p}$ , t.q. :

$$X = U \Sigma V^T,$$

avec  $\Sigma = \text{diag}(s_1, \dots, s_r)$ ,  $s_1 \geq \dots \geq s_r > 0$ , et  $r = \text{rang}(X)$ .

**Démonstration.** Diagonaliser  $X^T X$  [Golub and Van Loan, 2013]

· Les  $s_j \in \mathbb{R}^+$  sont les **valeurs singulières** de  $X$ , et les  $u_j$  (resp.  $v_j$ ) sont les **vecteurs singuliers** à gauche (resp. à droite)



# Décomposition en valeurs singulières (*singular value decomposition, SVD*)

· Si  $n < p$  :

$$X = U\Sigma V^T = [u_1, \dots, u_n] \left[ \begin{array}{ccc|c} s_1 & \cdots & 0 & \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & s_n & \end{array} \right]_{n \times (p-n)} \left[ \begin{array}{c} v_1^T \\ \vdots \\ v_p^T \end{array} \right] = \sum_{i=1}^n s_i u_i v_i^T$$

· Si  $p < n$  :

$$X = U\Sigma V^T = [u_1, \dots, u_n] \left[ \begin{array}{ccc|c} s_1 & \cdots & 0 & \\ \vdots & \ddots & \vdots & \\ 0 & \cdots & s_p & \\ \hline & & & 0_{(n-p) \times p} \end{array} \right] \left[ \begin{array}{c} v_1^T \\ \vdots \\ v_p^T \end{array} \right] = \sum_{i=1}^p s_i u_i v_i^T$$

· Si  $n < p$  :

$$X = U\Sigma V^T = [u_1, \dots, u_n] \begin{bmatrix} s_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & s_n \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_n^T \end{bmatrix} = \sum_{i=1}^n s_i u_i v_i^T$$

· Si  $p < n$  :

$$X = U\Sigma V^T = [u_1, \dots, u_p] \begin{bmatrix} s_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & s_p \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_p^T \end{bmatrix} = \sum_{i=1}^p s_i u_i v_i^T$$

**SVD réduite** : on ne garde que les éléments utiles avec  $r = \min(n, p)$  :

$$X = U\Sigma V^T = \sum_{i=1}^r s_i u_i v_i^T = U_r \Sigma_r V_r^T, \quad (13)$$

avec  $s_i > 0$  pour tout  $i \in \llbracket 1, r \rrbracket$ , et  $U_r = [u_1, \dots, u_r]$ ,  $V_r = [v_1, \dots, v_r]$

**SVD compacte** : on ne garde que les  $r = \text{rang}(X)$  valeurs singulières non-nulles

**Rem** : les  $u_i$  (resp. les  $v_i^T$ ) sont orthonormés et engendrent le même espace que celui engendré par les colonnes (resp. les lignes) de  $X$

$$\text{vect}(x_1, \dots, x_p) = \text{vect}(u_1, \dots, u_r)$$

## Rappel :

$$X = U\Sigma V^T \Rightarrow \Sigma = U^T X V$$

- La première valeur singulière  $s_1$  s'obtient du problème d'optimisation :

$$s_1 = \begin{cases} \max_{u \in \mathbb{R}^n, v \in \mathbb{R}^p} & u^T X v \\ \text{s.c.} & \|u\|^2 = 1, \\ & \|v\|^2 = 1 \end{cases}$$

- Lagrangien :

$$\mathcal{L}(u, v) = u^T X v - \lambda_1(\|u\|^2 - 1) - \lambda_2(\|v\|^2 - 1)$$

- Conditions nécessaires d'optimalité (CNO) :

$$\begin{cases} \nabla_u \mathcal{L} = Xv - 2\lambda_1 u = 0 \\ \nabla_v \mathcal{L} = X^T u - 2\lambda_2 v = 0 \end{cases} \iff \begin{cases} Xv = 2\lambda_1 u \\ X^T u = 2\lambda_2 v \end{cases} \Rightarrow \begin{cases} X^T X v = \alpha v \\ X X^T u = \alpha u \end{cases}$$

avec  $\alpha = 4\lambda_1\lambda_2$ , et donc  $v$  et  $u$  sont des vecteurs propres de  $X^T X$  et de  $X X^T$

- Considérons que  $X \in \mathbb{R}^{n \times p}$  admet pour SVD  $X = \sum_{i=1}^r s_i u_i v_i^\top$  avec  $r = \text{rang}(X)$
- La **pseudo-inverse** de  $X$ , dénotée par  $X^+ \in \mathbb{R}^{p \times n}$ , est définie par :

$$\begin{aligned} X^+ &:= (U \Sigma V^\top)^{-1} = (V^\top)^{-1} \Sigma^{-1} U^{-1} \\ &= V \text{diag}(1/s_1, \dots, 1/s_r) U^\top \\ &= \sum_{i=1}^r \frac{1}{s_i} v_i u_i^\top \end{aligned}$$

- Si  $X = \sum_{i=1}^n s_i u_i v_i^\top \in \mathbb{R}^{n \times n}$  est inversible alors  $X^+ = X^{-1}$

- Les fonctions SVD et pseudo-inverse sont disponibles dans les librairies numériques classiques, par exemple Numpy
- **SVD** :  $U, s, V = \text{np.linalg.svd}(X)$ 
  - $\Sigma = \text{diag}(s)$  ( $X = U \cdot \text{dot}(\text{np.diag}(s)) \cdot \text{dot}(V)$ )
  - Variantes compactes ou non par l'option : `full_matrices=True/False`
- **Pseudo-inverse** : `np.linalg.pinv(X)`

· Partons de la SVD  $X = U\Sigma V^T = \sum_{i=1}^r s_i u_i v_i^T$  :

$$\|X\beta - y\|^2 = \left\| U\Sigma V^T \beta - y \right\|^2$$

· Partons de la SVD  $X = U\Sigma V^T = \sum_{i=1}^r s_i u_i v_i^T$  :

$$\begin{aligned}\|X\beta - y\|^2 &= \left\| U\Sigma V^T \beta - y \right\|^2 = \left\| U\Sigma V^T - UU^T y \right\|^2 \\ &= \left\| \sum_{i=1}^r s_i u_i v_i^T \beta - \sum_{i=1}^n u_i u_i^T y \right\|^2 \\ &= \left\| \sum_{i=1}^r u_i (s_i v_i^T \beta - u_i^T y) - \sum_{i=r+1}^n u_i u_i^T y \right\|^2 \\ &= \left\| \sum_{i=1}^r u_i (s_i v_i^T \beta - u_i^T y) \right\|^2 + \left\| \sum_{i=r+1}^n u_i u_i^T y \right\|^2 \\ &= \sum_{i=1}^r \left( s_i v_i^T \beta - u_i^T y \right)^2 + \sum_{i=r+1}^n (u_i^T y)^2\end{aligned}$$



$$\|X\beta - y\|^2 = \sum_{i=1}^r (s_i v_i^\top \beta - u_i^\top y)^2 + \sum_{i=r+1}^n (u_i^\top y)^2 \geq \sum_{i=r+1}^n (u_i^\top y)^2$$

· L'égalité s'obtient si

$$\beta = \sum_{i=1}^r \frac{\langle y, u_i \rangle}{s_i} v_i = \sum_{i=1}^r \frac{v_i u_i^\top}{s_i} y,$$

$$\text{or } X^+ = \sum_{i=1}^r \frac{1}{s_i} v_i u_i^\top$$

· Ainsi une solution des moindres carrés peut s'écrire:

$$\hat{\beta} = X^+ y \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2} \|X\beta - y\|^2$$

**Rem :** l'ensemble de toutes les solutions est:

$$\left\{ X^+ y + \sum_{i=r+1}^p \alpha_i v_i, (\alpha_{r+1}, \dots, \alpha_p) \in \mathbb{R}^{p-r} \right\},$$

où  $X^+ y$  est la solution de norme  $\|\cdot\|$  minimale

- Sous l'hypothèse de bruit "blanc" (i.e.,  $\mathbb{E}(\varepsilon) = 0$ ):

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(X^+y) = \sum_{i=1}^r v_i v_i^T \beta^* = \Pi_l \beta^*$$

- $\Pi_l$  : projecteur sur l'espace des lignes de  $X$

$$\Pi_l = \sum_{i=1}^r v_i v_i^T = X^+ X$$

- $\Pi_c$  : projecteur sur l'espace des colonnes de  $X$

$$\Pi_c = \sum_{i=1}^r u_i u_i^T = X X^+$$

**Rem :** si  $r = \text{rang}(X) = n$  on retrouve que les conditions d'optimalité sont sans biais :  $\mathbb{E}(\hat{\beta}) = \beta^*$

- Sous l'hypothèse de bruit "blanc" (i.e.,  $\mathbb{E}(\varepsilon) = 0$ ):

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(X^+y) = \sum_{i=1}^r v_i v_i^\top \beta^* = \Pi_l \beta^*$$

- $\Pi_l$  : projecteur sur l'espace des lignes de  $X$

$$\Pi_l = \sum_{i=1}^r v_i v_i^\top = X^+ X$$

- $\Pi_c$  : projecteur sur l'espace des colonnes de  $X$

$$\Pi_c = \sum_{i=1}^r u_i u_i^\top = X X^+$$

**Rem :** si  $r = \text{rang}(X) = n$  on retrouve que les conditions d'optimalité sont sans biais :  $\mathbb{E}(\hat{\beta}) = \beta^*$

**Démonstration.**

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}(X^+[X\beta^* + \varepsilon]) = \sum_{i=1}^r \frac{1}{s_i} v_i u_i^\top \sum_{j=1}^r s_j u_j v_j^\top \beta^* = \sum_{i=1}^r v_i v_i^\top \beta^* = \Pi_l \beta^*$$

- Hypothèse : bruit blanc et homoscedastique (i.e.,  $\mathbb{E}(\varepsilon) = \mathbf{0}$ ,  $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$ )

$$\text{cov}(\hat{\beta}) = \text{cov}(X^+y) = \sum_{i=1}^r \frac{\sigma^2}{S_i^2} \mathbf{v}_i \mathbf{v}_i^\top$$

**Rem :** si  $\text{rang}(X) = n$ , on retrouve  $\text{cov}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$

- Hypothèse : bruit blanc et homoscedastique (i.e.,  $\mathbb{E}(\varepsilon) = \mathbf{0}$ ,  $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$ )

$$\text{cov}(\hat{\beta}) = \text{cov}(X^+y) = \sum_{i=1}^r \frac{\sigma^2}{S_i^2} v_i v_i^\top$$

**Rem :** si  $\text{rang}(X) = n$ , on retrouve  $\text{cov}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$

**Démonstration.**

$$\text{cov}(\hat{\beta}) = \text{cov}(X^+[X\beta^* + \varepsilon]) = \mathbb{E}(X^+ \varepsilon \varepsilon^\top (X^+)^{\top}) = X^+ \mathbb{E}(\varepsilon \varepsilon^\top) (X^+)^{\top} = \sum_{i=1}^r \frac{\sigma^2}{S_i^2} v_i v_i^\top$$

- Hypothèse : modèle homoscédastique (*i.e.*,  $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$ ) et que  $X$  est de plein rang

$$R_{pred}(y, \hat{y}) = \mathbb{E}\|X\beta^* - X\hat{\beta}\|_2^2 = \sigma^2 \text{rang}(X)$$

- Hypothèse : modèle homoscédastique (i.e.,  $\mathbb{E}(\varepsilon\varepsilon^\top) = \sigma^2 I_n$ ) et que  $X$  est de plein rang

$$R_{pred}(y, \hat{y}) = \mathbb{E} \|X\beta^* - X\hat{\beta}\|_2^2 = \sigma^2 \text{rang}(X)$$

## Démonstration.

$$\begin{aligned} R_{pred}(y, \hat{y}) &= \mathbb{E} \left( (\hat{\beta} - \beta^*)^\top X^\top X (\hat{\beta} - \beta^*) \right) \\ &= \mathbb{E} \left( (X^+ [X\beta^* + \varepsilon] - \beta^*)^\top X^\top X (X^+ [X\beta^* + \varepsilon] - \beta^*) \right) \\ &= \mathbb{E} \left( (X^+ \varepsilon + [\Pi_l - I_p] \beta^*)^\top X^\top X (X^+ \varepsilon + [\Pi_l - I_p] \beta^*) \right) \\ &= \mathbb{E} \left( (X^+ \varepsilon)^\top X^\top X (X^+ \varepsilon) \right) + \underbrace{(\beta^*)^\top (\Pi_l - I_p)^\top X^\top X (\Pi_l - I_p) \beta^*}_{=0} \\ &= \mathbb{E} \left( (X X^+ \varepsilon)^\top X X^+ \varepsilon \right) \\ &= \mathbb{E} \left( \text{tr}(X X^+ \varepsilon (X X^+ \varepsilon)^\top) \right) \quad (a^\top b = \text{tr}(ba^\top)) \\ &= \text{tr}(\Pi_c \mathbb{E}(\varepsilon\varepsilon^\top) \Pi_c^\top) \\ &= \sigma^2 \text{tr}(\Pi_c \Pi_c^\top) = \sigma^2 \text{rang}(\Pi_c) = \sigma^2 r = \sigma^2 \text{rang}(X) \end{aligned}$$

- Prenons  $\hat{\beta} = X^+y$  comme solution des moindres carrés.
- Supposons qu'on observe maintenant  $y + \Delta$  où  $\|\Delta\| \ll \|y\|$
- Alors l'estimateur des moindres carrés pour  $y + \Delta$  par  $X$  donne

$$\hat{\beta}^\Delta = X^+(y + \Delta)$$

$$\hat{\beta}^\Delta = \hat{\beta} + X^+ \Delta$$

$$\hat{\beta}^\Delta = \hat{\beta} + \sum_{i=1}^r \frac{1}{s_i} v_i u_i^\top \Delta$$

**Rem :** Noter l'influence des “petites” valeurs singulières



- Remèdes possibles contre les mauvais “conditionnements”
  - Régulariser le spectre / les valeurs singulières
  - Contraindre les coefficients de  $\hat{\beta}$  à n’être pas trop grands
- Une solution rendant ces deux points de vue équivalents : *Ridge Regression* / Régularisation de Tychonoff

## Réduction de la dimension

---

- Soit  $X = [x_1^\top, \dots, x_n^\top]^\top \in \mathbb{R}^{n \times p}$  la matrice des variables explicatives
- L'ACP (de niveau  $k \in \mathbb{N}$ ) consiste à effectuer la SVD de  $X$ , et à ne garder que les  $k$  axes principaux pour représenter le nuage

$$X = \sum_{i=1}^k s_i u_i v_i^\top$$

- On appelle **axes principaux** (ou **axes factoriels**) les vecteurs  $v_1, \dots, v_k \in \mathbb{R}^p$ . En général  $k \ll p$  (e.g.,  $k = 2$ , pour une visualisation plane)
- Les nouvelles variables  $c_i = Xv_i$  sont appelées **composantes**

**Rem :** On doit recentrer les points pour qu'ils aient une moyenne nulle  $X \leftarrow [(x_1 - \bar{x}_n)^\top, \dots, (x_n - \bar{x}_n)^\top]^\top = X - \mathbf{1}_n \bar{x}_n^\top$  (on peut aussi mettre à l'échelle pour avoir un écart-type similaire par caractéristique (*feature*) )

- On applique d'abord une ACP sur  $X$  afin d'obtenir les nouvelles variables

$$z_i = Xv_i,$$

pour tout  $i \in \{1, \dots, k\}$

- Ensuite, on considère le modèle linéaire sur  $([z_1, \dots, z_k] \in \mathbb{R}^{n \times k}, y \in \mathbb{R}^n)$  :

$$y = \bar{y}1_n + \sum_{i=1}^k \theta_i z_i$$

- L'estimateur  $\hat{\theta} \in \mathbb{R}^{k+1}$  s'obtient comme suit **[exercice]**

$$\hat{\theta}_i = \frac{\langle z_i, y \rangle}{\langle z_i, z_i \rangle}$$

- Car  $z_1, \dots, z_k$  s'obtiennent par des combinaisons linéaires de  $X$ , alors

$$\hat{\beta}_k = \sum_{i=1}^k \hat{\theta}_i v_i$$

**Rem :** Si  $k = p$ , on revient sur l'estimation des moindres carrés

# Régression des moindres carrés partiels (*Partial Least Squares, PLS*)

- PLS est une méthode itérative permettant la réduction de la dimension

---

## Algorithm Partial Least Squares [Hastie et al., 2009]

---

**Données d'entrée :**  $(X, y)$ , nombre d'axes PLS  $k \ll p$

1: Recentrer et standardiser  $X$

2: Définir  $\hat{y}^{(0)} = \bar{y}1_n$  et  $x_j^{(0)} = x_j$  pour tout  $j \in \{1, \dots, p\}$

3: **for**  $i = 1, \dots, k$  **do**

4:  $z_i = \sum_{j=1}^p \hat{\phi}_{i,j} x_j^{(i-1)}$  où  $\hat{\phi}_{i,j} = \langle x_j^{(i-1)}, y \rangle$

5:  $\hat{y}^{(i)} = \hat{y}^{(i-1)} + \hat{\theta}_i z_i$  avec  $\hat{\theta}_i = \frac{\langle z_i, y \rangle}{\langle z_i, z_i \rangle}$

6: Orthogonaliser chaque  $x_j^{(i-1)}$  par rapport à  $z_i$  :

$$x_j^{(i)} = x_j^{(i-1)} - \frac{\langle z_i, x_j^{(i-1)} \rangle}{\langle z_i, z_i \rangle} z_i$$

7: Car  $z_1, \dots, z_k$  sont linéaire par rapport à  $x_j$ , alors

$$\hat{y} = \hat{y}^{(p)} = \bar{y}1_n + \sum_{i=1}^k \hat{\theta}_i z_i$$

- 
- Si  $k = p$ , on trouve une solution équivalente à celle des moindres carrés

## PCA

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^p} \quad & \text{var}(X\alpha) \\ \text{s.c.} \quad & \|\alpha\|^2 = 1, \\ & \alpha^\top S v_\ell = 0, \quad \forall \ell = 1, \dots, i-1 \end{aligned}$$

## PLS

$$\begin{aligned} \max_{\alpha \in \mathbb{R}^p} \quad & \text{corr}^2(y, X\alpha) \text{var}(X\alpha) \\ \text{s.c.} \quad & \|\alpha\|^2 = 1, \\ & \alpha^\top S v_\ell = 0, \quad \forall \ell = 1, \dots, i-1 \end{aligned}$$

- $S$  est la matrice de covariance empirique de  $x_j$
- Les conditions  $\alpha^\top S v_\ell = 0$  rassurent que  $z_i = X\alpha$  n'est pas corrélée avec les combinaisons linéaires précédentes  $z_\ell = Xv_\ell$

- G. H. Golub and C. F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 4 edition, 2013.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer Series in Statistics. Springer, New York, second edition, 2009.
- Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct):2825–2830, 2011.