

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import matplotlib.colors as colors
```

Détection d'anomalies

Mise en situation

Je travaille pour une banque. Je dispose des données de transactions bancaires par client, et j'aimerais détecter une utilisation suspecte - une anomalie.

Je ne sais pas à quoi ça peut ressembler une anomalie : je ne peux pas vraiment faire de classification supervisée. Et puis quand bien même : un nouveau type de fraude pourrait apparaître, que je ne saurais pas reconnaître...

Détection d'anomalies par analyse des distributions gaussiennes

Il existe plusieurs classes d'algorithmes pour la détection d'anomalies. Celle présentée ici n'est pas la meilleure mais c'est probablement la plus simple à comprendre / appliquer.

Accessoirement, elle n'est pas mauvaise non plus...

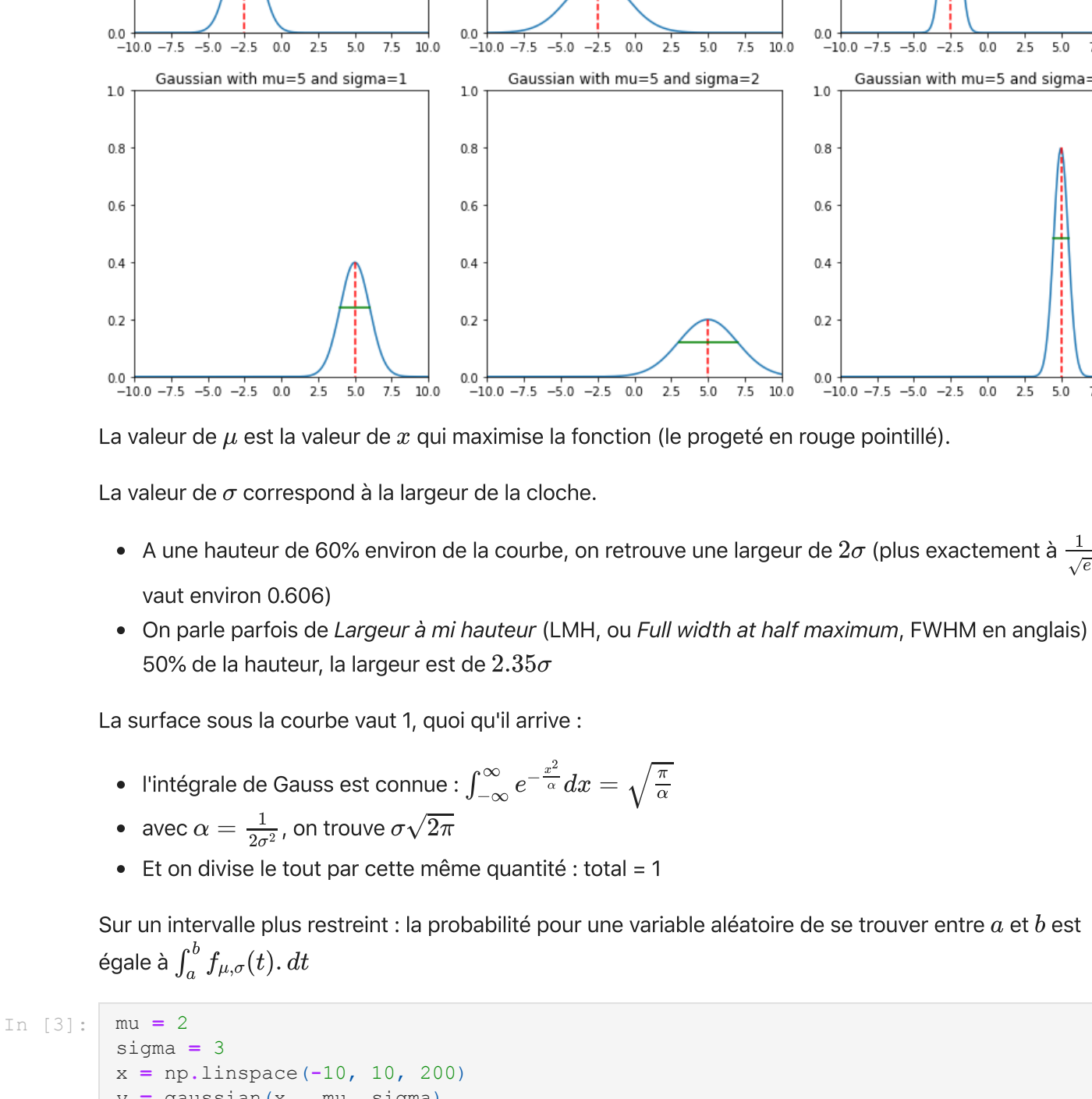
Distribution gaussienne ?

Une variable aléatoire X est gaussienne (ou normale) lorsqu'elle suit une *loi normale*, dont la densité de probabilité est :

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

avec μ la moyenne des valeurs de cette variable, et σ l'écart-type. On note la distribution $\mathcal{N}(\mu, \sigma^2)$, le carré de l'écart type étant la *variance*.

Comme je suis lucide, je me doute que c'est par forcément plus clair pour autant, alors un peu de visuel ne fera pas de mal : voici quelques courbes représentatives de f pour différentes valeurs de μ et de σ



La valeur de μ est la valeur de x qui maximise la fonction (le projeté en rouge pointillé).

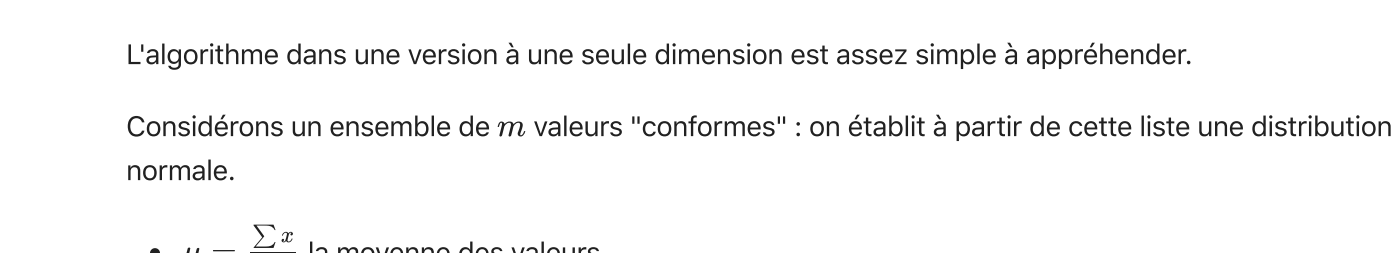
La valeur de σ correspond à la largeur de la cloche.

- A une hauteur de 60% environ de la courbe, on retrouve une largeur de 2σ (plus exactement à $\frac{1}{\sqrt{e}}$ qui vaut environ 0.606)
- On parle parfois de *Largeur à mi hauteur* (LMH, ou *Full width at half maximum*, FWHM en anglais) : à 50% de la hauteur, la largeur est de 2.35σ

La surface sous la courbe vaut 1, quoi qu'il arrive :

- l'intégrale de Gauss est connue : $\int_{-\infty}^{\infty} e^{-\frac{x^2}{2\sigma^2}} dx = \sqrt{\pi}$
- avec $\alpha = \frac{1}{2\sigma^2}$, on trouve $\sigma\sqrt{2\pi}$
- Et on divise le tout par cette même quantité : total = 1

Sur un intervalle plus restreint : la probabilité pour une variable aléatoire de se trouver entre α et β est égale à $\int_{\alpha}^{\beta} f_{\mu,\sigma}(t).dt$



On peut voir sur cette figure que la probabilité d'avoir une valeur de 0 est plus élevée que d'avoir une valeur de -5, mais moins que d'avoir 2.3.

Pour être tout à fait exact, la probabilité d'avoir exactement 0 est nulle, comme toute probabilité d'avoir une valeur exacte sur ce genre de distributions. Mais la probabilité d'avoir une valeur dans un rayon ϵ autour de 0 est plus importante que pour le même rayon autour de -5 mais moins que pour 2.3.

Détection d'anomalies à une dimension

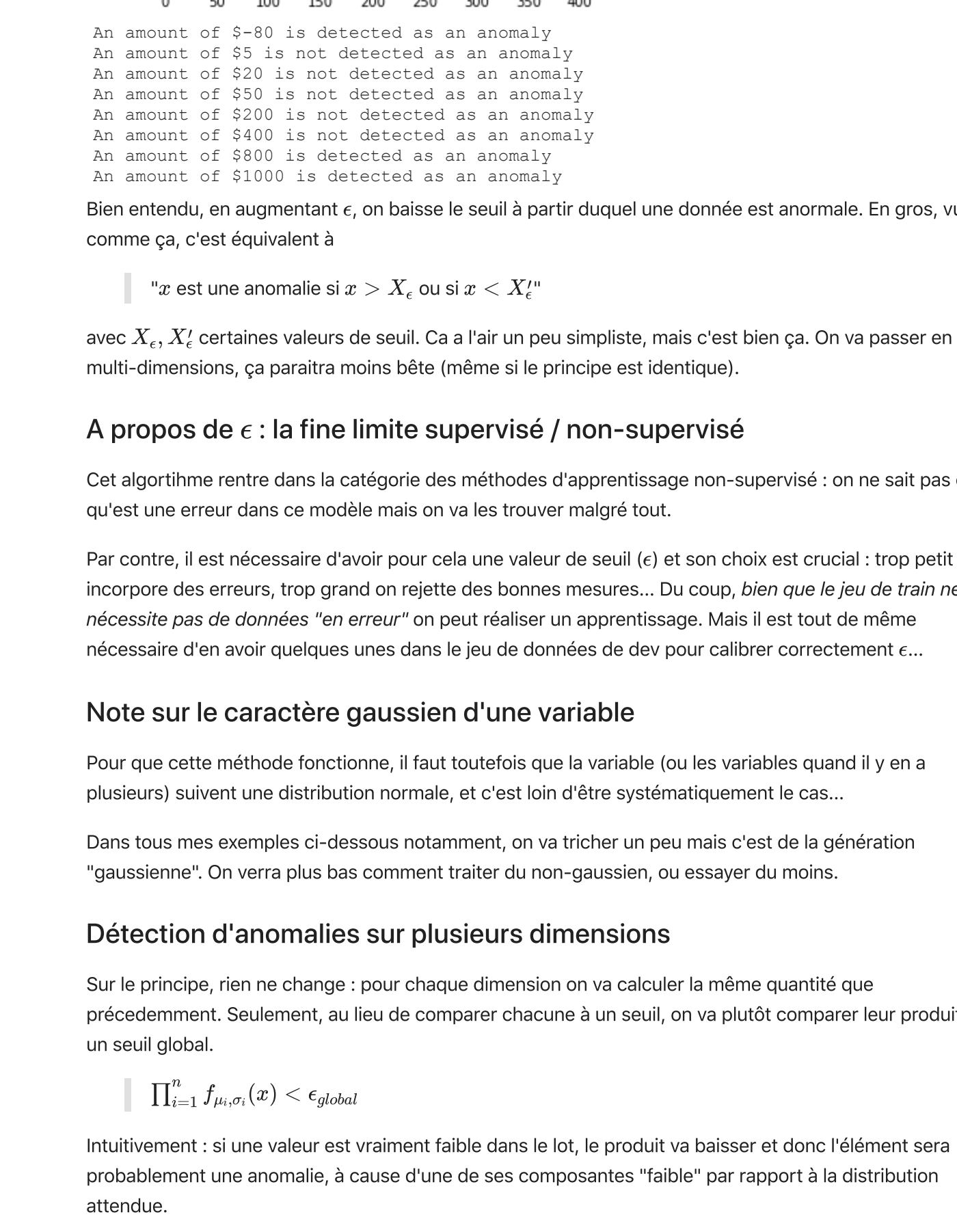
L'algorithmie dans une version à une seule dimension est assez simple à appréhender.

Considérons un ensemble de m valeurs "conformes" : on établit à partir de cette liste une distribution normale.

- $\mu = \frac{\sum x_i}{m}$ la moyenne des valeurs
- $\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{m}}$ l'écart type

Pour toute nouvelle valeur x , on considérera qu'elle est une anomalie si $f_{\mu,\sigma}(x) < \epsilon$, avec ϵ une constante, très petite, à définir.

Exemple :



An amount of 5-80 is detected as an anomaly
An amount of 85 is not detected as an anomaly
An amount of \$200 is not detected as an anomaly
An amount of \$50 is not detected as an anomaly
An amount of \$200 is not detected as an anomaly
An amount of \$400 is not detected as an anomaly
An amount of \$800 is detected as an anomaly
An amount of \$1000 is detected as an anomaly

Bien entendu, en augmentant ϵ , on baisse le seuil à partir duquel une donnée est anormale. En gros, vu comme ça, c'est équilibré à

$$\mu \pm \epsilon \text{ "x est une anomalie si } x > X_\epsilon \text{ ou si } x < X_\epsilon \text{"}$$

avec X_ϵ, X_ϵ certaines valeurs de seuil. Ca a l'air un peu simpliste, mais c'est bien ça. On va passer en multi-dimensions, ça paraîtra moins bête (même si le principe est identique).

A propos de ϵ : la fine limite supervisé / non-supervisé

Cet algorithme rentre dans la catégorie des méthodes d'apprentissage non-supervisé : on ne sait pas ce qu'est une erreur dans ce modèle mais on va la trouver malgré tout.

Par contre, il est nécessaire d'avoir pour cela une valeur de seuil (ϵ) et son choix est crucial : trop petit on incorpore des erreurs, trop grand on rejette des bonnes mesures... Du coup, *bien que le jeu de train ne nécessite pas de données "en erreur"* on peut réaliser un apprentissage. Mais il est tout de même nécessaire d'en avoir quelques unes dans le jeu de données de dev pour calibrer correctement ϵ ...

Note sur le caractère gaussien d'une variable

Pour que cette méthode fonctionne, il faut toutefois que la variable (ou les variables quand il y en a plusieurs) suivent une distribution normale, et c'est loin d'être systématiquement le cas...

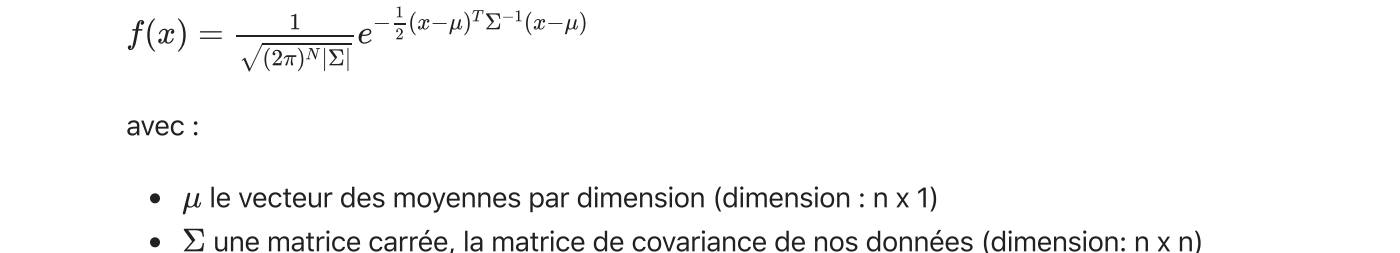
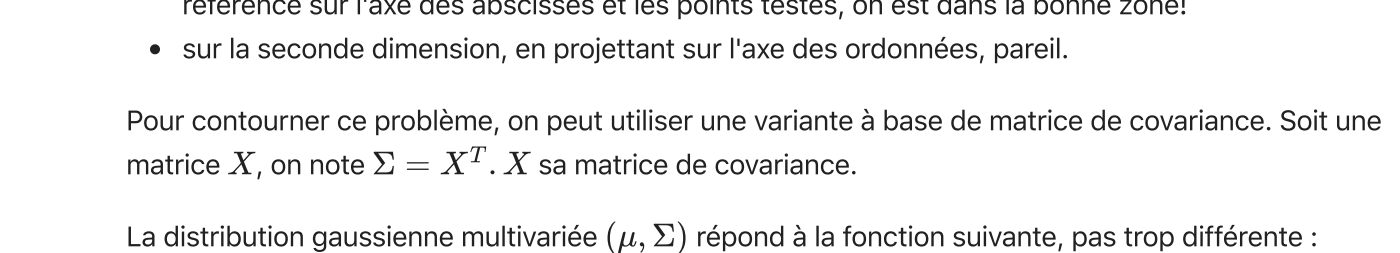
Dans tous mes exemples ci-dessous notamment, on va tricher un peu mais c'est de la génération "gaussienne". On verra plus bas comment traiter du non-gaussien, ou essayer du moins.

Détection d'anomalies sur plusieurs dimensions

Sur le principe, rien ne change : pour chaque dimension on va calculer la même quantité que précédemment. Seulement, au lieu de comparer chacune à un seuil, on va plutôt comparer leur produit à un seuil global.

$$\prod_{i=1}^n f_{\mu_i, \sigma_i}(x_i) < \epsilon_{global}$$

Intuitivement : si une valeur est vraiment faible dans le lot, le produit va baisser et donc l'élément sera probablement une anomalie, à cause d'une de ses composantes "faible" par rapport à la distribution attendue.



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

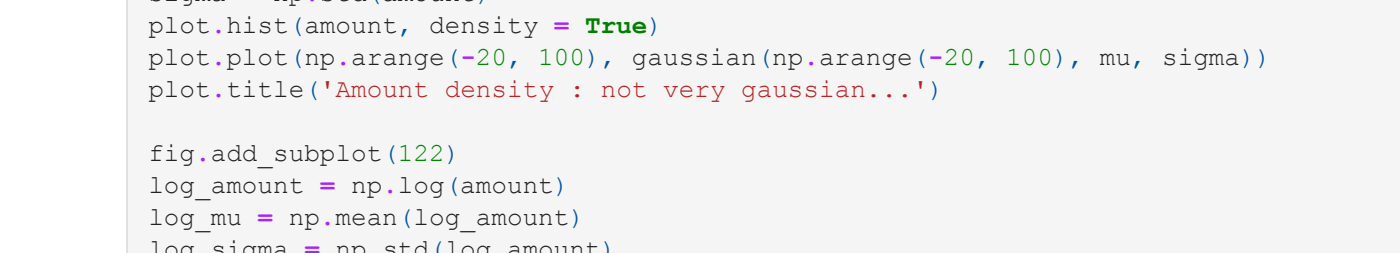
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

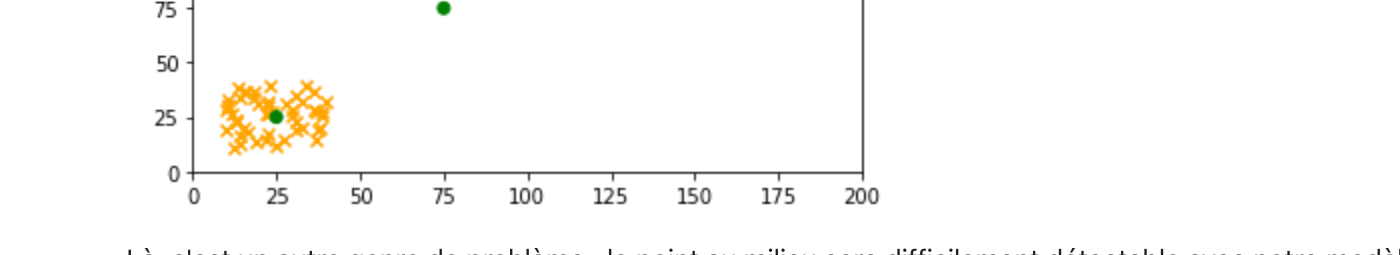
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

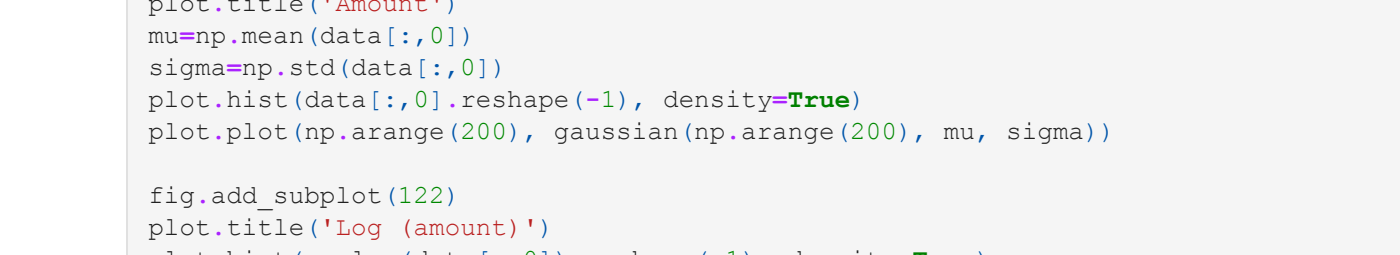
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

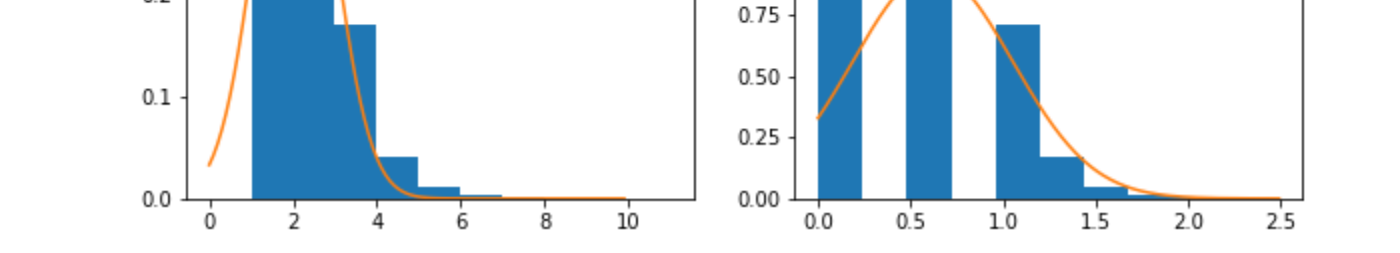
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

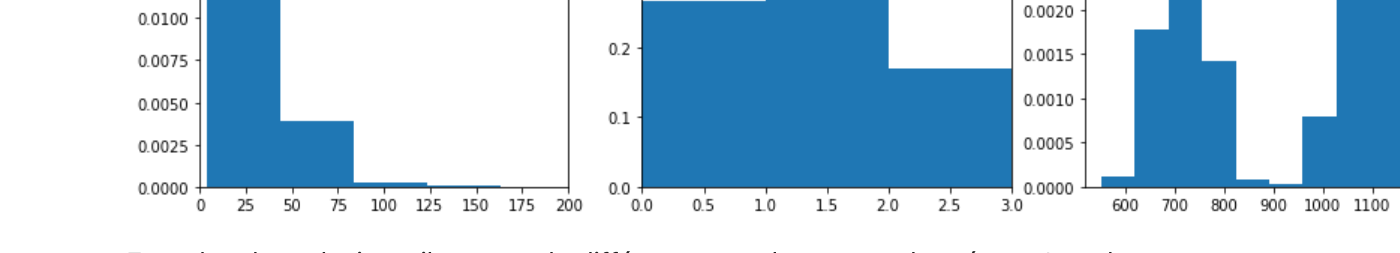
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l'une augmente, l'autre aussi.

Et on a clairement un souci : certains points sont verts, ils devraient être rouges, et le seul rouge devrait être vert. On a mis ϵ à un pour un milliard, augmenter cette valeur exclura les 2 mauvais points verts, mais va exclure encore plus fortement le mauvais point rouge. Inversement, baisser cette valeur passera tout en vert.

Explication :

- sur la première dimension, on est entre 200 et 1200 à chaque fois : si on projette les données de référence sur l'axe des abscisses et les points testés, on est dans la bonne zone!
- sur la seconde dimension, en projetant sur l'axe des ordonnées, pareil.

Pour contourner ce problème, on peut utiliser une variante à base de matrice de covariance. Soit une matrice Σ , on note $\Sigma = X^T \cdot X$ sa matrice de covariance.

La distribution gaussienne multivariée (μ, Σ) répond à la fonction suivante, pas trop différente :

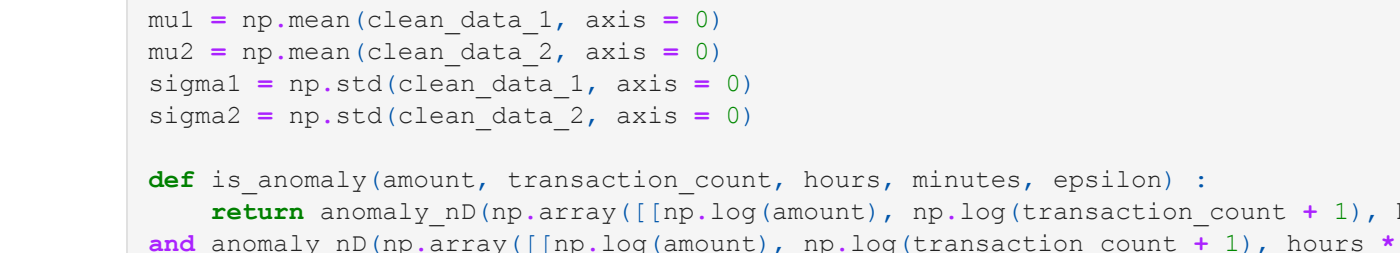
$$f(x) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

avec :

- μ le vecteur des moyennes par dimension (dimension : $n \times 1$)
- Σ une matrice carrée, la matrice de covariance de nos données (dimension : $n \times n$)
- x la multivarié que l'on teste (dimension $n \times 1$)
- $|\Sigma|$ le déterminant de la matrice Σ

Dans notre implémentation

- le dataset comporte un point par colonne, μ est calculé comme un vecteur ligne, les x aussi : on permute les transpositions $\rightarrow (x - \mu)^T \Sigma^{-1} (x - \mu)^T$
- le résultat pour m tests simultanés est une matrice carrée au lieu d'un vecteur : on en prend la diagonale



On voit bien sur le graphique que les deux mesures sont un peu dépendantes : quand l