



**POLYTECHNIQUE
MONTREAL**

LE GÉNIE
EN PREMIÈRE CLASSE

INF8460 –Introduction au traitement de la langue naturelle

Automne 2019

TP No. [1]

Groupe [2]

[1853357] – [El-Asmar Jad]

[1847744] – [Elakhrass Mohammed Ali]

[1847193] – [Hasrouni Elio]

Soumis à : Félix Martel

[24 septembre 2019]

Table des matières

1.1 Exploration des données.....	3
Les statistiques descriptives sur Shakespeare	3
1.2 Modèle de langue n-gramme.....	3
Les K n-grams les plus fréquents ($K=20$, $n=1,2,3$)	3
1.5 Évaluation des modèles	4
Les perplexités des modèles	4
Les graphes reportant la perplexité en fonction du paramètre gamma	5
1.6 Génération de texte	6
Les séquences générées sur le corpus Trump	6

1.1 Exploration des données

Les statistiques descriptives sur Shakespeare

-- shakespeare_train --

Nombre de tokens: 1046885

Nombre de types: 52482

Les 15 mots les plus fréquents du vocabulaire:

[(',', 79977), (':', 44576), ('.', 33719), ('the', 22775), ('I', 21423), ('and', 16464), (';', 15364), ('to', 15150), ('of', 14698), ('you', 12259), ('a', 11966), ('my', 10474), ('?', 10167), ('in', 9286), ('is', 9020)]

Ratio: 19.94750581151633

Nombre de lemmes distincts: 52482

Nombre de racines (stems) distinctes: 52482

-- shakespeare_test --

Nombre de tokens: 1991

Nombre de types: 46

Les 15 mots les plus fréquents du vocabulaire:

[(',', 164), ('the', 40), ('thou', 40), ('to', 31), ('thy', 31), ('in', 30), (':', 28), ('of', 28), ('s', 27), ('and', 26), ('.', 26), (';', 26), ('that', 21), ('And', 21), ('thee', 20)]

Ratio: 43.28260869565217

Nombre de lemmes distincts: 46

Nombre de racines (stems) distinctes: 46

1.2 Modèle de langue n-gramme

Les K n-grams les plus fréquents (K=20, n =1,2,3)

Pour n=1

[(',', 79977), ((), ':', 44576), ((), '.', 33719), ((), 'the', 22775), ((), 'I', 21423), ((), 'and', 16464), ((), ';', 15364), ((), 'to', 15150), ((), 'of', 14698), ((), 'you', 12259), ((), 'a', 11966), ((), 'my', 10474), ((), '?', 10167), ((), 'in', 9286), ((), 'is', 9020), ((), '!', 8765), ((), 'not', 8579), ((), 'that', 8062), ((), 'me', 7539), ((), 'And', 7057)]

Pour n=2

[('(', '</s>', 33401), (('?', '</s>', 10067), (('!', '</s>', 8479), ((';', 'and', 4873), ((';', 'I', 4411), ((':', 'I', 4097), ((';', 'And', 2852), ((';', 'my', 2107), (('!',

'am', 1806), (('<s>', 'I', 1781), (('!', '"', 1715), ((';', 'and', 1673), (('you', ' ', 1626), ((';', 'sir', 1623), (('!', 'have', 1575), (('!', 'will', 1532), ((';', 'that', 1510), (('in', 'the', 1467), ((';', 'And', 1380), (('to', 'the', 1331)]

Pour n=3

[('(', '</s>', '</s>', 33401), (('?', '</s>', '</s>', 10067), (('!', '</s>', '</s>', 8479), (('<s>', '<s>', 'I', 1781), (('<s>', '<s>', 'KING', 1308), ((';', 'my', 'lord', 956), ((';', 'sir', ' ', 867), (('<s>', '<s>', 'First', 773), (('<s>', '<s>', 'What', 744), (('<s>', '<s>', 'O', 733), (('you', ' ', '</s>', 692), (('it', ' ', '</s>', 676), (('<s>', '<s>', 'But', 662), (('me', ' ', '</s>', 618), (('him', ' ', '</s>', 588), ((':', 'O', ' ', 586), ((';', 'Why', ' ', 569), ((';', 'Ay', ' ', 558), (('"', '</s>', '</s>', 534), (('<s>', '<s>', 'The', 522)]

1.5 Évaluation des modèles

Les perplexités des modèles

n= 1

Modèle MLE: 821.9874763092969

Modèle Laplace: 821.3779146370026

Modèle Lidstone :

	Gamma									
	1e-05	3.59381366 38046256e -05	0.00012915 496650148 84	0.000464 15888336 12782	0.00166810 053720005 92	0.005994842 503189409	0.02154 4346900 318846	0.07742636 826811278	0.27825594 02207126	1.0
Perplexité	821.987456 5068733	821.987405 1446898	821.987220 5776148	821.9865 57523999 7	821.984177 8106788	821.9756665 422401	821.945 6049392 913	821.844233 6975195	821.560334 9352118	821.377914 6370026

n= 2

Modèle MLE: inf

Modèle Laplace: 1916.377487869097

Modèle Lidstone :

	Gamma									
	1e-05	3.59381366 38046256e -05	0.00012915 496650148 84	0.000464 15888336 12782	0.00166810 053720005 92	0.005994842 503189409	0.02154 4346900 318846	0.07742636 826811278	0.27825594 02207126	1.0
Perplexité	3374.32987 72060334	2344.31257 9149452	1654.19386 87797224	1209.770 33368285 47	943.737771 9562788	811.7734409 186937	795.270 4134676 615	909.782864 2953434	1227.71516 21028695	1916.37748 7869097

n= 3

Modèle MLE: inf

Modèle Laplace: 8695.41062097361

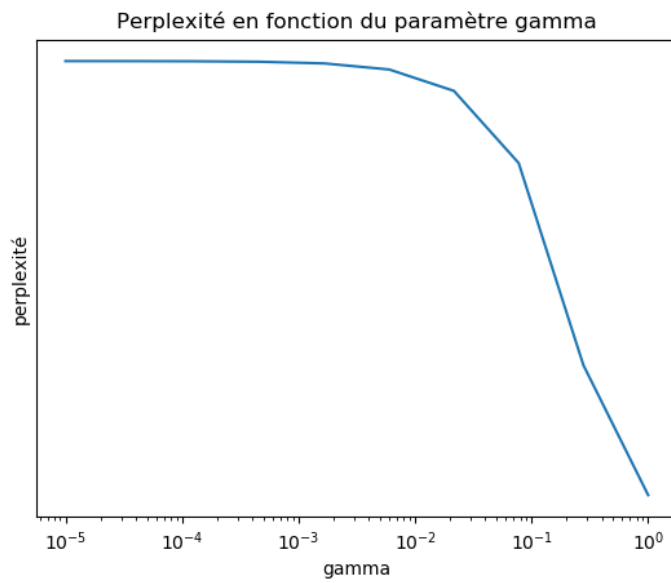
Modèle Lidstone :

	Gamma									
	1e-05	3.59381366 38046256e -05	0.00012915 496650148 84	0.000464 15888336 12782	0.00166810 053720005 92	0.005994842 503189409	0.02154 4346900 318846	0.07742636 826811278	0.27825594 02207126	1.0
Perplexité	17475.6008 90810958	10530.4336 09227226	6821.38258 6957344	4959.136 04279484 3	4147.33521 2902398	4004.801572 927881	4400.33 3522675 265	5326.44350 0861947	6792.79905 5926109	8695.41062 097361

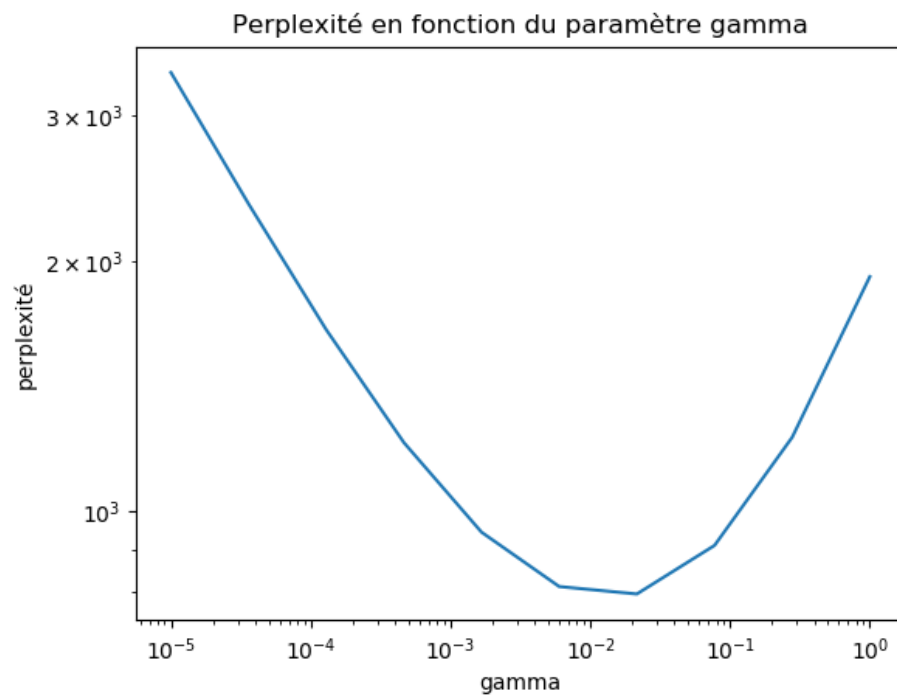
Le meilleur modèle est celui de Laplace puisque la perplexité est la plus petite pour les trois valeurs de n. En ce qui concerne le modèle MLE, pour n = 1, il s'agit de la formule suivante qui s'applique pour les probabilités $P(W_i) = \frac{COUNT(W_i)}{N}$. En effet, la variable N n'est jamais nulle puisqu'il s'agit du nombre de mots (tokens). Cependant pour n=2 et n=3 il s'agit de la formule suivante qui s'applique $P(w_i | w_{i-1}) = \frac{count(w_{i-1}, w_i)}{count(w_{i-1})}$. Effectivement, pour cette formule si dans le corpus de test nous rencontrons un mot qui n'est pas présent dans le corpus d'entraînement, count(w_{i-1}) aura une valeur de 0. Cela provoque une division par 0 donnant qui tend vers l'infini.

Les graphes reportant la perplexité en fonction du paramètre gamma

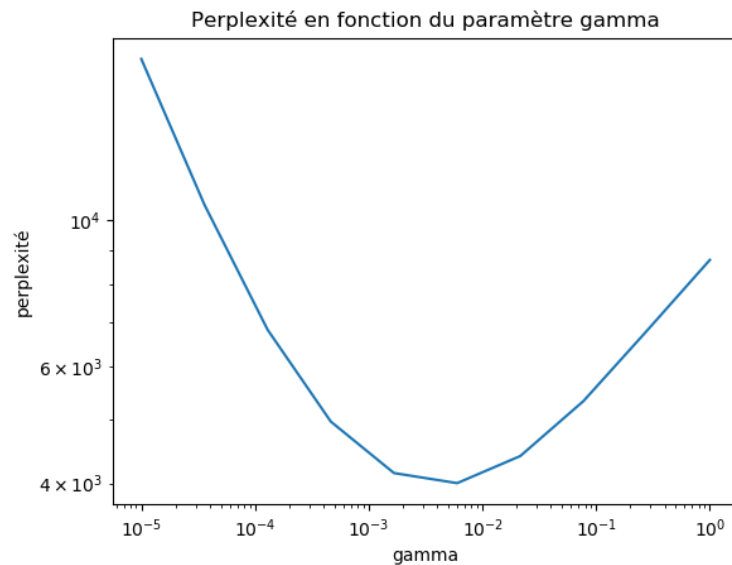
Pour n=1 :



Pour n=2 :



Pour $n=3$:



Nous remarquons que quand γ tend vers zéro, la valeur de la perplexité est la plus grande. Pour $n = 1$, le gamma avec la valeur de 1.0 donne le meilleur résultat, soit de 821.3779146370026. Pour $n = 2$, il s'agit du gamma avec la valeur 0.021544346900318846 donnant la meilleure perplexité de 795.2704134676615. Pour $n = 3$, il s'agit du gamma avec la valeur de 0.005994842503189409 donnant la meilleure perplexité de 4004.801572927881.

1.6 Génération de texte

Les séquences générées sur le corpus Trump

n= 1

['a', 'was', '.', ',', 'lemons', 'Wasserman', 'DonaldWillBeTheGOPNominee', 'everyone', '.', 'Are', 'Supreme', 'you', 'URL', 'end', ':', 'all', ',', 'with', 'Museum', 'you']

['Indiana', 'basket', ',', ',', 'nytimes', '"', 'Debates20_', 'on', 'vote', '#border', 'it', 'just', 'TRUMP', 'night', ':', 'but', ':', '!']

n= 2

['#Texas', ',', 'we', 'can', 'make', 'it', 'up', 'in', 'this', 'year', 'America', 'For', 'reasons', 'only', 'had', 'some', 'of', 'the', 'family']

['Details', 'to', 'being', 'Native', 'American', 'Senator', 'Lindsey', 'Graham', '@realDonaldTrump', 'Love', 'U', 'REALLY', 'loves', '@club4growth', 'has', 'predicted', 'Trumps', 'Speech']

n= 3

['Such', 'wonderful', 'people', 'living', 'in', 'poverty', ',', 'education', 'and', 'safety', 'of', 'Americans', 'than',
'offending', 'people', '""', 'Tonight', 'I', 'will', 'fix']

[''', '@_HankRearden', ':', 'I', 'endorse', '@realDonaldTrump', 'URL', '""', '""', '@JulesSiscoe', ':', 'We',
'need', 'to', 'MAKE', 'AMERICA', 'GREAT']