Para que los modelos funcionen óptimamente, se procede a realizar la estandarización del corpus proporcionado. El primer elemento detectado es que el corpus utiliza marcas para identificar patrones sintácticos. Por ejemplo -Fpa-, -Fpt- y 0, pero como el objetivo del laboratorio es generar texto fluido con n-gramas, se procede a filtrarlos. In [4]: def removeSpecialTags(sentences): Elimina tokens como -Fpa-, -Fpt- (rodeados por guiones) y tokens como \*0\*, \*1\* (placeholders de anotación sintáctica) en el corpus cess\_esp. cleaned = [] for sentence in sentences: cleaned sentence = [ word **for** word **in** sentence if not (word.startswith('-') and word.endswith('-')) # -Fpa-, -Fptand not re.fullmatch(r"\\*\d+\\*", word) # \*0\*, \*1\*, etc. cleaned.append(cleaned sentence) return cleaned comments = removeSpecialTags(corpus) In [5]: comments = [' '.join(sent) for sent in comments] comments[:10] Out[5]: ['El grupo estatal Electricité\_de\_France EDF anunció hoy , jueves , la compra del 51\_por\_ciento de la empresa mexicana Electricidad\_Águila\_de\_Altamira EAA , creada por el japonés Mitsubishi\_Corporation para poner\_en\_marcha una central de gas de 495 megavatios .', 'Una portavoz de EDF explicó a EFE que el proyecto para la construcción de Altamira 2 , al norte de Tampico , prevé la utilización de gas natural como combustible principal en una central de ciclo combinado que debe empezar a funcionar en mayo \_del\_2002 .', 'La electricidad producida pasará a la red eléctrica pública de México en virtud de un acuerdo de venta de energía de EAA con la Comisión Federal de Electricidad CFE por una duración de 25 años .', 'EDF , que no quiso revelar cuánto pagó por su participación mayoritaria en EAA , intervendrá como asistente en la construcción de Altamira\_2 y , posteriormente , se encargará de explotarla como principal accionista .', 'EDF y Mitsubishi participaron en 1998 en la licitación de licencias para construir centrales eléctricas en México y se quedaron con dos cada una : Río\_Bravo y Saltillo para la compañía francesa y Altamira y Tuxpán para la japonesa .', 'EDF tiene previsto invertir 194 millones de euros 186 millones de dólares en la central de Río Bravo, con una potencia de 495 megavatios, y 134 millones de euros 28 millones de dólares en Saltillo, que como la primera funcionará con gas na tural y cuya potencia prevista es de 247 megavatios .', 'La alcaldesa de Málaga y cabeza de lista del PP al Congreso por esta provincia , Celia\_Villalobos , pidió hoy a los militantes de esta formación que sepan " administrar la victoria " , porque " no vale la revancha , el insulto o el ataque , e so es para ellos " .', 'En una intervención ante militantes populares en el hotel de la capital malagueña elegido como sede de la noche electoral del PP, Villalobos dio las gracias " a los militantes de muchos años que hoy tienen una emoción especial ", y se confe só también " especialmente emocionada " en sus quintas elecciones generales .', 'La diputada electa transmitió " un abrazo y mi cariño a José María Aznar " , a quien definió como " el hombre que ha sabido llevar al PP hasta la victoria " , y dio las gracias a la " gente que ha sabido ilusionar y convencer a muchos para qu e apuesten " por este partido .', 'A todas estas personas " no las vamos a defraudar " , aseguró Villalobos , que subrayó que Aznar " seguirá gobernando desde el diálogo , sin prepotencia , con honradez y preocupado por los problemas , y seguiremos siendo los mismos " .'] Se buscar analizar el efecto de la estandarización para identificar si existen apariciones de nuevas palabras con más repeticiones o simplemente, aumenta el número de las que actualmente se encuentran como más frecuentes. Para esto, se mostraran las 8 palabras más frecuentes antes del proceso y después del mismo. In [6]: comments = [comment.strip() for comment in comments if comment.strip()] def top\_n\_words(comments, n=8): from collections import Counter words = [word for comment in comments for word in comment.split()] return Counter(words).most\_common(n) print("Top 8 common words:") print(top\_n\_words(comments)) Top 8 common words: [(',', 11420), ('de', 10234), ('la', 6412), ('.', 5866), ('que', 5552), ('el', 5199), ('en', 4340), ('y', 4235)] Para verificar si será necesario aplicar la transformación de mayusculas a mínusculas, se desarrolla una función que permite observar si una palabra tiene una equivalente en con algún caracter en mayúscula o minúscula en el corpus. In [7]: def encontrar\_variaciones\_de\_casing(comments, max\_resultados=3): formas\_por\_palabra = {} resultados\_mostrados = 0 for comment in comments: for word in comment.split(): base = word.lower() if base not in formas\_por\_palabra: formas\_por\_palabra[base] = set() formas\_por\_palabra[base].add(word) for base, formas in formas\_por\_palabra.items(): if len(formas) > 1: print(f"Variantes de '{base}': {formas}") resultados\_mostrados += 1 if resultados\_mostrados >= max\_resultados: break print("Words with case variations:") encontrar\_variaciones\_de\_casing(comments) Words with case variations: Variantes de 'el': {'EL', 'el', 'El'} Variantes de 'grupo': {'Grupo', 'grupo'} Variantes de 'anunció': {'anunció', 'Anunció'} Observamos que sí existen casos con variantes mayusculas y minusculas, por lo que sí se considera necesaria su aplicación. Además, se procede a analizar si ocurre lo mismo pero con signos de puntuación. In [8]: def encontrar\_variaciones\_de\_puntuacion(comments, max\_resultados=3): formas\_por\_palabra = {} resultados\_mostrados = 0 signos = ".!?,;:()\"'¿;" for comment in comments: for word in comment.split(): base = word.strip(signos) if base not in formas\_por\_palabra: formas\_por\_palabra[base] = set() formas\_por\_palabra[base].add(word) for base, formas in formas\_por\_palabra.items(): if len(formas) > 1: print(f"Variantes de puntuación para '{base}': {formas}") resultados\_mostrados += 1 if resultados\_mostrados >= max\_resultados: break print("Words with punctuation variations:") encontrar\_variaciones\_de\_puntuacion(comments) Words with punctuation variations: Variantes de puntuación para '': {',', ';', '!', '¡', "'", '¿', '"', ':', '?', '...'} Variantes de puntuación para 'PP': {'PP', 'PP.'} Variantes de puntuación para '2': {'2', '2.'} De esta forma identificamos que también se producen ocurrencias, por lo que se aplica el procedimiento de eliminación de puntuación In [9]: def erasePunctuation(comments): return [re.sub(r'[^\w\s]', '', line) for line in comments] Se procede con la estandarización y se evalua el resultado. In [10]: comments = [comment.lower() for comment in comments] comments = erasePunctuation(comments) print("Top 8 common words after standardization:") print(top\_n\_words(comments)) Top 8 common words after standardization: [('de', 10286), ('la', 6925), ('el', 6013), ('que', 5570), ('en', 4643), ('y', 4349), ('los', 3189), ('a', 3020)] [('de', 10286), ('la', 6925), ('el', 6013), ('que', 5570), ('en', 4643), ('y', 4349), ('los', 3189), ('a', 3020)] Se puede apreciar que existe un aumento en la frecuencia de algunas palabras, por lo tanto, la estandarización es efectiva. In [11]: comments[:10] Out[11]: ['el grupo estatal electricité\_de\_france edf anunció hoy jueves la compra del 51\_por\_ciento de la empresa mexicana electricidad\_águila\_de\_altamira eaa creada por el japonés mitsubishi\_corporation para poner\_en\_marcha una central de gas de 49 5 megavatios ', 'una portavoz de edf explicó a efe que el proyecto para la construcción de altamira\_2 al norte de tampico prevé la utilización de gas natural como combustible principal en una central de ciclo combinado que debe empezar a funcionar en mayo\_d 'la electricidad producida pasará a la red eléctrica pública de méxico en\_virtud\_de un acuerdo de venta de energía de eaa con la comisión\_federal\_de\_electricidad cfe por una duración de 25 años ', 'edf que no quiso revelar cuánto pagó por su participación mayoritaria en eaa intervendrá como asistente en la construcción de altamira 2 y posteriormente se encargará de explotarla como principal accionista ', 'edf y mitsubishi participaron en 1998 en la licitación de licencias para construir centrales eléctricas en méxico y se quedaron con dos cada una río bravo y saltillo para la compañía francesa y altamira y tuxpán para la japonesa ', 'edf tiene previsto invertir 194 millones de euros 186 millones de dólares en la central de río bravo con una potencia de 495 megavatios y 134 millones de euros 28 millones de dólares en saltillo que como la primera funcionará con gas natur al y cuya potencia prevista es de 247 megavatios ', 'la alcaldesa de málaga y cabeza de lista del pp al congreso por esta provincia celia\_villalobos pidió hoy a los militantes de esta formación que sepan administrar la victoria porque no vale la revancha el insulto o el ataque eso es pa ra ellos ', 'en una intervención ante militantes populares en el hotel de la capital malagueña elegido como sede de la noche electoral del pp villalobos dio las gracias a los militantes de muchos años que hoy tienen una emoción especial y se confesó t ambién especialmente emocionada en sus quintas elecciones generales ', 'la diputada electa transmitió un abrazo y mi cariño a josé maría aznar a quien definió como el hombre que ha sabido llevar al pp hasta la victoria y dio las gracias a la gente que ha sabido ilusionar y convencer a muchos para que apues ten por este partido ', 'a todas estas personas no las vamos a defraudar aseguró villalobos que subrayó que aznar seguirá gobernando desde el diálogo sin prepotencia con honradez y preocupado por los problemas y seguiremos siendo los mismos '] Finalmente se identifica que existen secuencias de 2 o más espacios en blanco en la separación de algunas palabras, para evitar obtener bigramas incorrectos o que en algun punto de la estandarización afecten, se procede a normalizarlos. In [12]: def normalizeSpaces(lines): return [re.sub(r'\s+', ' ', line).strip() for line in lines] In [13]: comments = normalizeSpaces(comments) comments[:10] Out[13]: ['el grupo estatal electricité\_de\_france edf anunció hoy jueves la compra del 51\_por\_ciento de la empresa mexicana electricidad\_águila\_de\_altamira eaa creada por el japonés mitsubishi\_corporation para poner\_en\_marcha una central de gas de 495 m 'una portavoz de edf explicó a efe que el proyecto para la construcción de altamira 2 al norte de tampico prevé la utilización de gas natural como combustible principal en una central de ciclo combinado que debe empezar a funcionar en mayo del \_2002', 'la electricidad producida pasará a la red eléctrica pública de méxico en virtud de un acuerdo de venta de energía de eaa con la comisión federal de electricidad cfe por una duración de 25 años', 'edf que no quiso revelar cuánto pagó por su participación mayoritaria en eaa intervendrá como asistente en la construcción de altamira 2 y posteriormente se encargará de explotarla como principal accionista', 'edf y mitsubishi participaron en 1998 en la licitación de licencias para construir centrales eléctricas en méxico y se quedaron con dos cada una río\_bravo y saltillo para la compañía francesa y altamira y tuxpán para la japonesa', 'edf tiene previsto invertir 194 millones de euros 186 millones de dólares en la central de río\_bravo con una potencia de 495 megavatios y 134 millones de euros 28 millones de dólares en saltillo que como la primera funcionará con gas natural y cuya potencia prevista es de 247 megavatios', 'la alcaldesa de málaga y cabeza de lista del pp al congreso por esta provincia celia villalobos pidió hoy a los militantes de esta formación que sepan administrar la victoria porque no vale la revancha el insulto o el ataque eso es para ello s', 'en una intervención ante militantes populares en el hotel de la capital malagueña elegido como sede de la noche electoral del pp villalobos dio las gracias a los militantes de muchos años que hoy tienen una emoción especial y se confesó tambi én especialmente emocionada en sus quintas elecciones generales', 'la diputada electa transmitió un abrazo y mi cariño a josé maría aznar a quien definió como el hombre que ha sabido llevar al pp hasta la victoria y dio las gracias a la gente que ha sabido ilusionar y convencer a muchos para que apuesten por este partido', 'a todas estas personas no las vamos a defraudar aseguró villalobos que subrayó que aznar seguirá gobernando desde el diálogo sin prepotencia con honradez y preocupado por los problemas y seguiremos siendo los mismos'] De esta forma se procede a almacenar el corpus estandarizado. In [14]: with open("corpus.txt", "w", encoding="utf-8") as file: for comment in comments: file.write(comment + "\n") Construcción de modelos de n-gramas Los modelos basados en bigramas requieren de una tokenización de los elementos para luego ser procesados mediante la formula condicional P(wn | wn-1, ..., wn-k+1) usando frecuencias relativas. Con ayuda de nitk se construye un modelo de conteo condicional que permite calcular por cada palabra wn, cuántas veces aparece wn-1 después de la misma. Lo cual permite realizar el calculo de las probabilidades condicionales contando el total de ocurrencias por cada wn, usandolo como denominador de la frecuencia mencionada anteriormente. Básicamente, se construyen los modelos en base a la siguiente fórmula:  $[P(w_n \in w_{n-1}, \cdot w_{n-k+1}) = \frac{C(w_{n-k+1}, \cdot w_n)}{C(w_{n-k+1}, \cdot w_{n-1})}]$ donde (k = 1, 2, 3)In [15]: # Carga y tokenización del corpus estandarizado with open("corpus.txt", "r", encoding="utf-8") as file: lines = file.readlines() tokens = [] for line in lines: words = line.strip().split() tokens.extend(words) tokens[:10] Out[15]: ['el', 'grupo', 'estatal', 'electricité\_de\_france', 'edf', 'anunció', 'hoy', 'jueves', 'la', 'compra'] In [16]: def trainNgramModel(comments, n): tokens = [] for sentence in comments: tokens.extend(sentence.strip().split()) ngram\_list = list(ngrams(tokens, n)) cfd = ConditionalFreqDist() for ng in ngram\_list: context = () if n == 1 else tuple(ng[:-1]) word = ng[0] if n == 1 else ng[-1]cfd[context][word] += 1  $model = \{\}$ for context in cfd: model[context] = {

Aunque en el ejercicio anterior se construyeron modelos de n-gramas sobre todo el corpus, para este análisis de espacidad es necesario separar una parte del corpus como conjunto de prueba. Esto permite evaluar cuántos n-gramas nuevos aparecen que no fueron observados en

Se infiere con esto que a mayor valor de n, mayor espacidad: el modelo necesita más combinaciones de palabras. De esta forma se limita la capacidad del modelo de generar texto o calcular la probabilidad para frases que son nuevas.

Antes de realizar la implementación, los modelos se entrenarán en este caso con todo el corpus ya que en incisos posteriores lo que se busca evaluar es la perplejidad por lo que no existe la necesidad de hacer la separación en train y test.

UVG

UNIVERSIDAD

**DEL VALLE** 

CAMPUS SUR

In [1]: import nltk

Out[1]: True

[nltk\_data]

nltk.download('cess esp')

[nltk data] Downloading package cess esp to

[nltk\_data] Package cess\_esp is already up-to-date!

In [2]: from nltk.lm.preprocessing import padded everygram pipeline

Preprocesamiento del corpus

'\_\_total\_\_': cfd[context].N()

bigramsProbability.append(bigramProbability)

bigramsProbability.append(bigramProbability)

tokens.extend(sentence.strip().split())

context = () if n == 1 else tuple(ng[:-1])

if context not in model or word not in model[context]:

sparsity = calculateSparsity(model, test\_comments, n\_size)

print(f"Espacidad para n={n\_size}: {sparsity:.2f}% de n-gramas no vistos")

C(contexto, w) + 1

C(contexto) + V

In [24]: def interpolatedProbability(uni\_model, bi\_model, tri\_model, context, word, vocab\_size, lambdas=(0.4, 0.4, 0.2)):

test\_ngrams = getNgrams(test\_comments, n)

word = ng[0] if n == 1 else ng[-1]

model[context][word] = cfd[context][word]

bigramProbability, vocabSizeC = trainNgramModel(comments, n\_size)

bigramProbability, vocabSizeC = trainNgramModel(train\_comments, n\_size)

los datos de entrenamiento, lo cual es esencial para entender las limitaciones de los modelos. La separación se realizará con 75% entrenamiento y 25% prueba.

El suavizado en modelos de n-gramas es una técnica que ajusta las probabilidades para que incluso los n-gramas no observados tengan una probabilidad distinta de cero.

for word in cfd[context]:

vocab = set(tokens)

In [17]: bigramsProbability = []
 n = [1, 2, 3]
 for n\_size in n:

In [19]: bigramsProbability = []
for n size in n:

In [20]: def getNgrams(comments, n):
 tokens = []

missing = 0

In [22]: for i, n\_size in enumerate(n):

Suavizado de laplace

Interpolación lineal

unigram = ()

**Kneser-Ney** 

In [25]: def trainKneserNeyModel(tokens, n):

return model

bigramsVocabSize = []

import math

if N < n:</pre>

import math

**if** N < 3:

# Tabla vacía

"n": [],

tabla\_perplejidad = {

"Laplace": [],

"Kneser-Ney": []

# === Laplace ===

**if** n size == 3:

**if** n\_size **==** 1:

# Mostrar la tabla

**0** 1 497.11

**1** 2 3454.97

**2** 3 14720.46

display(df\_perplejidad)

else:

"Interpolación": [],

for i, n size in enumerate(n):

N = len(tokens)

 $log_prob_sum = 0.0$ 

N = len(tokens)

 $log_prob_sum = 0.0$ 

In [26]: bigramsProbability = []

n = [1, 2, 3]
for n\_size in n:

 $P_{\text{Laplace}}(w \mid \text{contexto}) =$ 

return model, len(vocab)

Análisis de espacidad

In [18]: splitIndex = int(len(comments) \* 0.75)

train\_comments = comments[:splitIndex]
test\_comments = comments[splitIndex:]

for sentence in comments:

total = len(test\_ngrams)

for ng in test\_ngrams:

missing += 1

return (missing / total) \* 100

model = bigramsProbability[i]

Espacidad para n=1: 16.26% de n-gramas no vistos Espacidad para n=2: 62.92% de n-gramas no vistos Espacidad para n=3: 92.60% de n-gramas no vistos

Implementación de suavizado

In [23]: def laplaceSmoothedProbability(model, context, word, vocab\_size):

total = model.get(context, {}).get('\_\_total\_\_', 0)

count = model.get(context, {}).get(word, 0)

 $P(w_n) = \lambda_1 P_{ ext{trigrama}} + \lambda_2 P_{ ext{bigrama}} + \lambda_3 P_{ ext{unigrama}}$ 

bigram = (context[-1],) if len(context) >= 1 else ()

p1 = laplaceSmoothedProbability(uni\_model, unigram, word, vocab\_size)
p2 = laplaceSmoothedProbability(bi\_model, bigram, word, vocab\_size)
p3 = laplaceSmoothedProbability(tri model, trigram, word, vocab size)

trigram = context if len(context) >= 2 else ()

return lambda1 \* p1 + lambda2 \* p2 + lambda3 \* p3

train\_data, vocab = padded\_everygram\_pipeline(n, [tokens])

bigramProbability, vocabSizeC = trainNgramModel(comments, n\_size)

lambda3, lambda2, lambda1 = lambdas

model = KneserNeyInterpolated(order=n)

bigramsProbability.append(bigramProbability)

In [28]: def calculatePerplexityLaplace(tokens, model, vocab\_size, n):

In [27]: kneser\_models = [trainKneserNeyModel(tokens, n) for n in [1, 2, 3]]

bigramsVocabSize.append(vocabSizeC)

Cálculo de perplejidad

return float('inf')

for i in range(N - n + 1):

return float('inf')

for i in range(N - 3 + 1):

context = ngram[:-1]

word = ngram[-1]

ngram = tuple(tokens[i:i + n])

log\_prob\_sum += math.log(prob)

ngram = tuple(tokens[i:i + 3])

log\_prob\_sum += math.log(prob)

sentence\_tokens = sentence.strip().split()

tabla\_perplejidad["n"].append(n\_size)

model\_laplace = bigramsProbability[i]
vocab\_laplace = bigramsVocabSize[i]

# === Interpolación (solo para n = 3) ===

sentence\_tokens,

vocab\_laplace

bigramsProbability[0],
bigramsProbability[1],
bigramsProbability[2],

tabla\_perplejidad["Laplace"].append(round(ppl\_laplace, 2))

tabla\_perplejidad["Interpolación"].append(round(ppl\_interp, 2))

ppl\_interp = calculatePerplexityInterpolated(

tabla\_perplejidad["Interpolación"].append("-")

tabla\_perplejidad["Kneser-Ney"].append("-")

df\_perplejidad = pd.DataFrame(tabla\_perplejidad)

691.0

n Laplace Interpolación Kneser-Ney

Generación de texto

for \_ in range(length):

for word in vocab:

return ' '.join(generated)

for \_ in range(length):

for word in vocab:

context,

len(vocab)

probs[word] = prob

generated.append(next\_word)

In [33]: def generateTextKneserNey(model, n, seed, length=15):

generated = seed.strip().split()

word,

# Selección ponderada

return ' '.join(generated)

In [34]: seed = "el"

length = 15

all\_tokens = []

# Laplace

))

--- n = 1 ---

--- n = 2 ---

Kneser-Ney:

--- n = 3 ---

Interpolación:

Kneser-Ney:

Discusión final

Impacto del suavizado en la perplejidad

calidad de los resultados en la generación de texto.

Coherencia del texto generado según n y el suavizado

¿Cómo influye la espacidad en el rendimiento del modelo?

Laplace:

Laplace:

Laplace:

# Kneser-Ney

**if** n\_size >= 2:

print("Laplace:")

**if** n\_size == 3:

for sentence in comments:

vocabulario = set(all\_tokens)

for i, n size in enumerate(n):

 $print(f"\n--- n = \{n\_size\} ---")$ 

# Interpolación (solo para n=3)

vocabulario,

model = kneser\_models[i]
print("\nKneser-Ney:")

print(generateTextKneserNey(model, n\_size, seed, length))

el intentona lo la iñaki\_saez pasará henry\_ford la de del quede de de días las la

el 12 y 18 sobre pezzo pero en las que los aniones o los numerosos desplazamientos

considerablemente la perplejidad y reflejando un mejor rendimiento del modelo en cuestión.

el comunicado oficial manuel\_lapuente ha sido arrasadas en el lehendakari de hielo y hora de ti

n\_size,
seed,
length

print("\nInterpolación:")

print(generateTextInterpolated(
 bigramsProbability[0],
 bigramsProbability[1],
 bigramsProbability[2],

probs = {}

 $probs = \{\}$ 

generated = seed.strip().split()

probs[word] = prob

generated.append(next\_word)

generated = seed.strip().split()

# === Kneser-Ney usando model.perplexity() de nltk.lm ===

test\_ngrams = list(ngrams(sentence\_tokens, n\_size))
ppl\_kn = kneser\_models[i].perplexity(test\_ngrams)

307.78

670.7

context = tuple(generated[-(n-1):]) if n > 1 else ()

In [32]: def generateTextInterpolated(unigrams, bigrams, trigrams, vocab, n, seed, length):

# Asegura que el contexto tiene la longitud adecuada

prob = interpolatedProbability(

unigrams, bigrams, trigrams,

context = tuple(generated[-(n-1):]) if n > 1 else ()

return ' '.join(generated + generated\_words)

all\_tokens.extend(sentence.strip().split())

prob = laplaceSmoothedProbability(model, context, word, vocab\_size)

next\_word = random.choices(list(probs.keys()), weights=probs.values())[0]

context = tuple(generated[-(n-1):]) if len(generated) >= n-1 else tuple(generated)

next\_word = random.choices(list(probs.keys()), weights=probs.values())[0]

# Usamos directamente .generate con num\_words=length y el context inicial

print(generateTextLaplace(bigramsProbability[i], vocabulario, bigramsVocabSize[i], n\_size, seed, length))

el imperceptible hinchados preservación heroína columna sugerencias vítores estéril clausurar bajar pasiones confiere asiente desconfianza masas

el mariano denodado recambio desprendimientos pierde incluyó acusada irak seguidos doctor fornidos la matriculados porque licencia

el bancarios mirarle ortega gabinete\_del\_gobierno turba protagonizaran concedió embargo opositores afectaron fernando\_morientes pablo\_pinillos apreciando opuesto autopista\_del\_oeste

Los resultados evidencian cómo las distintas técnicas de suavizado afectan significativamente la perplejidad del modelo. El suavizado de Laplace genera perplejidades excesivamente altas a medida que aumenta el valor de n, especialmente en n=3 donde supera los 14,000. Esto

A medida que aumenta n, los modelos pueden considerar más contexto anterior, lo cual mejora la fluidez y coherencia del texto generado. Con n=1 (Laplace), el texto resulta desorganizado y con poca estructura. Para n=2, los modelos empiezan a generar frases parcialmente

La espacidad tiene un impacto grande en el rendimiento de los modelos. Cuando el corpus no contiene suficientes combinaciones posibles de palabras, los modelos sin un buen suavizado (como Laplace) asignan malas probabilidades, elevando la perplejidad y afectando la

generación de texto. Modelos como Kneser-Ney son menos vulnerables a este problema, ya que estiman la probabilidad de un n-grama basándose también en cómo aparecen las palabras en otros contextos, permitiendo una mejor generalización incluso con datos escasos.

coherentes. Finalmente, con n=3, los textos generados mediante Interpolación y Kneser-Ney presentan oraciones más naturales y con una mejor conexión entre palabras. Esto sugiere que la capacidad de modelar el contexto largo, combinada con un suavizado eficiente, mejora la

indica que este método no maneja adecuadamente la esparsidad en n-gramas a medida que n es mayor. En cambio, otras técnicas como la Interpolación con 691.0 de perplejidad y Kneser-Ney de670.7 logran una mejor distribución de las probabilidades, reduciendo

generated\_words = list(model.generate(length, text\_seed=context))

# Unimos el seed original con las nuevas palabras generadas

In [31]: def generateTextLaplace(model, vocab, vocab size, n, seed, length=15):

tabla\_perplejidad["Kneser-Ney"].append(round(ppl\_kn, 2))

return math.exp(-log\_prob\_sum / (N - 3 + 1))

In [30]: sentence = "el proyecto de la construcción es grande y complejo"

return math.exp(-log\_prob\_sum / (N - n + 1))

context = () if n == 1 else tuple(ngram[:-1])

prob = laplaceSmoothedProbability(model, context, word, vocab\_size)

In [29]: def calculatePerplexityInterpolated(tokens, uni model, bi model, tri model, vocab size, lambdas=(0.2, 0.3, 0.5)):

prob = interpolatedProbability(uni\_model, bi\_model, tri\_model, context, word, vocab\_size, lambdas)

ppl\_laplace = calculatePerplexityLaplace(sentence\_tokens, model\_laplace, vocab\_laplace, n\_size)

word = ngram[0] if n == 1 else ngram[-1]

model.fit(train data, vocab)

return (count + 1) / (total + vocab\_size)

return list(ngrams(tokens, n))

In [21]: def calculateSparsity(model, test\_comments, n):

from nltk import ngrams, ConditionalFreqDist

from nltk.lm import KneserNeyInterpolated

from nltk.corpus import cess\_esp

from nltk.util import ngrams

corpus = cess\_esp.sents()

import pandas as pd

In [3]: # Load the CESS-ESP corpus

import random

import re

C:\Users\daher\AppData\Roaming\nltk data...

Universidad del Valle de Guatemala - UVG

Laboratorio 4: Modelos de Lenguaje con N-gramas, Espacidad y Suavizado

Facultad de Ingeniería - Computación Curso: Procesamiento de Lenguaje Natural Sección: 10

• Diego Alexander Hernández Silvestre - **21270** 

**Autor:**