

---

---

# Proyecto Final

# Entrega Final

Sandra Niño

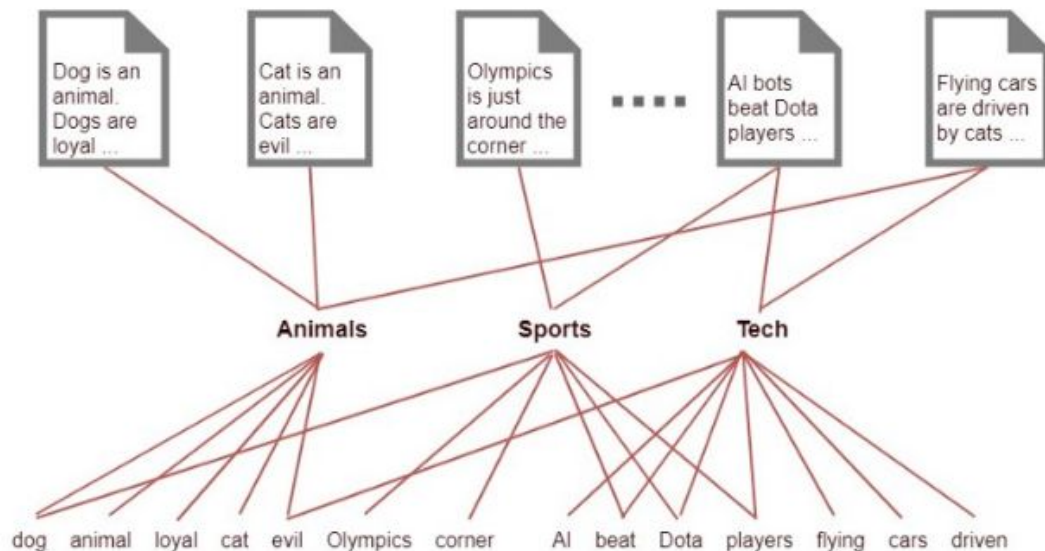
Johnatan Garzón

---

---

# LDA (Latent Dirichlet Allocation)

Cada documento puede ser descrito por una distribución de “topics” y cada topic puede ser descrito por una distribución de palabras.



# Procedimiento

Unión de  
datos en una  
sola columna



Se remueve  
valores null



Se remueve  
puntuación



Se pone en  
minúscula



Tokenización



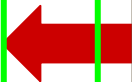
Se remueve  
stop-words



Bigramas y  
Trigramas



Lematización



LDA

# Remover Stop words

## Remover Stop words

Removemos las palabras que son conocidas como "Stop words": lo, la, con, antes, sin, de, entre otras, de los datos tokenizados.

```
stop_words_exceptions=["no"]
def remove_stopwords(texts):
    return [[word for word in simple_preprocess(str(doc)) if (word not in stop_words) or (word in stop_words_exceptions)] for doc in texts]

data_tokenized_nostops= remove_stopwords(data_tokenized)
```

---

{ 'demas\_personas', 'buena\_convivencia', 'sentirse\_tranquilo', 'cada\_persona', 'conmigo\_misma', 'poder\_compartir', 'respeto\_tolerancia', 'seres\_queridos', 'gracias\_dios', 'buenas\_personas', 'buena\_comunicacion', 'ninos\_jugando', 'bien\_conmigo', 'equidad\_social', 'conmigo\_mismo', 'mas\_alla', 'sigo\_mismo', 'poder\_convivir', 'entorno\_protector', 'ninos\_puedan', 'derechos\_humanos', 'iniciativas\_comunitarias', 'escuchar\_musica', 'buenas\_acciones', 'cultura\_ciudadana', 'puedan\_salir', 'mas\_importante', 'sana\_convivencia', 'poder\_salir', 'vivir\_tranquilo', 'buenos\_valores', 'no\_discriminacion', 'ser\_solidario', 'cualquier\_lugar', 'personas\_puedan', 'ser\_libre', 'ambiente\_sano', 'ninos\_ninas', 'tener\_buena', 'mismos\_derechos', 'ser\_humano', 'asi\_mismo', 'migo\_mismo', 'poder\_respirar', 'libre\_feliz', 'respeto\_hacia', 'libre\_expresion', 'buena\_relacion', 'convivencia\_pacifica', 'derechos\_sociales', 'si\_generamos', 'consigo\_mismo', 'justicia\_social', 'medio\_ambiente', 'seres\_vivos', 'fronteras\_pasar', 'aprender\_convivir', 'paz', 'ser\_tolerante', 'no\_violencia', 'trabajo\_empleo', 'seres\_humanos', 'recursos\_naturales', 'hacer\_actividades', 'jugar\_futbol' }

tamaño lista de bigramas sin stopwords: 1124

tamaño lista de bigramas sin stopwords únicos: 65

{ 'libre\_feliz', 'misma\_comunidad', 'ser\_humano', 'economica\_social', 'derechos\_sociales', 'mas\_alla', 'ser\_tolerantes', 'poder\_compartir', 'buenas\_personas', 'union\_familiar', 'ser\_libre', 'migo\_mismo', 'respeto\_hacia', 'sana\_convivencia', 'entorno\_protector', 'podamos\_convivir', 'asi\_mismo', 'tolerancia\_respeto', 'vida\_espiritual', 'poder\_hacer', 'cultura\_ciudadana', 'gracias\_dios', 'no\_pelear', 'poder\_disfrutar', 'están\_peleando', 'recursos\_naturales', 'respeto\_tolerancia', 'hacen\_parte', 'personas\_puedan', 'cada\_persona', 'siempre\_pensando', 'demas\_personas', 'conmigo\_misma', 'paz', 'familia\_compartir', 'sigo\_mismo', 'mas\_importante', 'no\_discriminacion\_derechos\_sociales', 'ninos\_ninas', 'derechos\_fundamentales', 'ala\_vida', 'servir\_ala', 'cada\_vez', 'fronteras\_pasar', 'derechos\_humanos', 'dia\_dia', 'no\_discriminacion', 'no\_haber', 'puedan\_salir', 'ninos\_jugando', 'no\_violencia\_derechos\_sociales', 'no\_violencia', 'seres\_vivos', 'bien\_conmigo', 'conceptos\_opiniones', 'buena\_relacion', 'hace\_feliz', 'seres\_queridos', 'poder\_respirar', 'ser\_solidario', 'ninos\_puedan', 'buena\_convivencia', 'conmigo\_mismo', 'aprender\_vivir', 'confianza\_interpersonal', 'buenas\_acciones', 'no\_discriminacion\_no\_violencia', 'escuchar\_musica', 'ser\_amable', 'sentirse\_tranquilo', 'vivir\_tranquilo', 'buena\_comunicacion', 'mismos\_derechos', 'ser\_feliz', 'consigo\_mismo', 'poder\_dialogar', 'ambiente\_sano', 'brindar\_amor', 'seres\_humanos', 'tener\_buena', 'jugar\_futbol', 'cualquier\_lugar', 'si\_generamos', 'aprender\_convivir', 'trabajo\_empleo', 'buenos\_valores', 'iniciativas\_comunitarias', 'libre\_expresion', 'poder\_convivir', 'poder\_jugar', 'participacion\_ciudadana', 'medio\_ambiente', 'poder\_salir', 'tener\_buena\_convivencia', 'tanta\_violencia', 'no\_aparece', 'salud\_trabajo', 'justicia\_social', 'deporte\_arte', 'escucho\_musica', 'buen\_trato', 'hacer\_actividades', 'convivencia\_pacifica', 'alegria\_felicidad', 'ser\_tolerante', 'equidad\_social' }

tamaño lista de trigramas sin stopwords: 1385

tamaño lista de trigramas sin stopwords únicos: 106

# Lematización

```
allowed_postags=['ADJ', 'VERB', 'ADV']
stopwordsToken=["demás", "tener", "mismo", "poder", "cada", "también", "hacer"]
unifiedWord={"respetar":"respeto", "violencia":"no_violencia", "musiciar":"musica", "guerra":"no_guerra"}
data_list=[]

for row in df.Narrativa_bigrams:
    token_list=""
    text=nlp(row)
    for token in text:
        if (token.pos_ in allowed_postags) and (token.lemma_ not in stopwordsToken) and (token.text not in stopwordsToken):
            if token.lemma_ in unifiedWord:
                token_list+=unifiedWord[token.lemma_]+" "
            else :
                token_list+=token.lemma_+" "

        elif (token.is_stop is not True) and (token.lemma_ not in stopwordsToken) and (token.text not in stopwordsToken):
            if token.text in unifiedWord:
                token_list+=unifiedWord[token.text]+" "
            else :
                token_list+=token.text+" "
    data_list.append(token_list.rstrip())
```



## Resultado de limpieza de datos



# LDA

Se crea una matriz Documento-Palabra como input para LDA

```
vectorizer = CountVectorizer(analyzer='word',  
                             min_df=10,          # minimo numero de ocurrencias  
                             token_pattern='[a-z\_A-Z]{3,}', # num chars > 3, que contenga bigramas (palabra_palabra)  
                             )
```

```
data_vectorized = vectorizer.fit_transform(df['Narrativa_lemmatized'])
```

```
df.Narrativa_lemmatized[0]
```

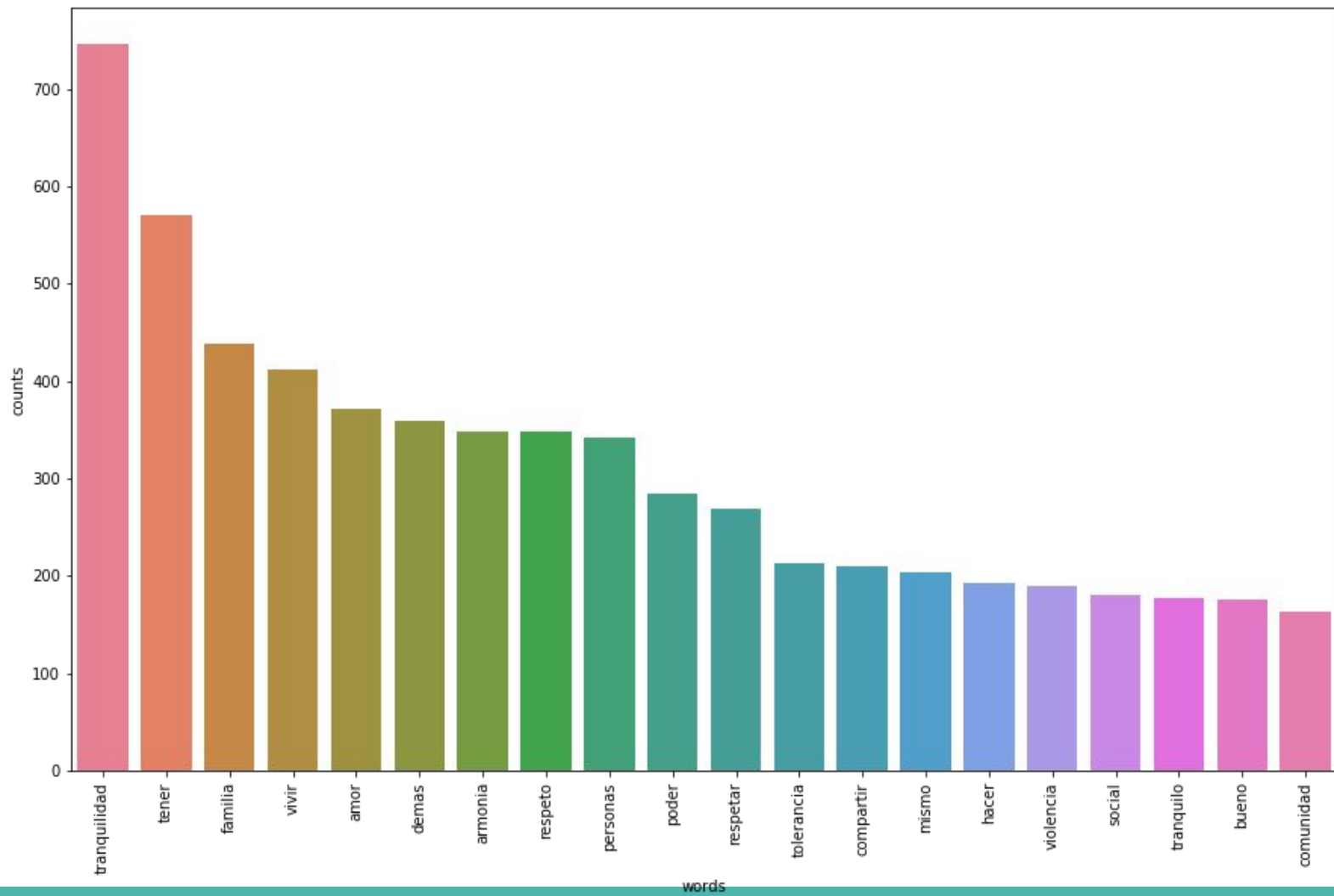
```
'tranquilidad convivencia casa educacion colectividad respeto tranquilidad acceso educacion humanizacion espacio siem  
pre camino'
```

```
print(data_vectorized)
```

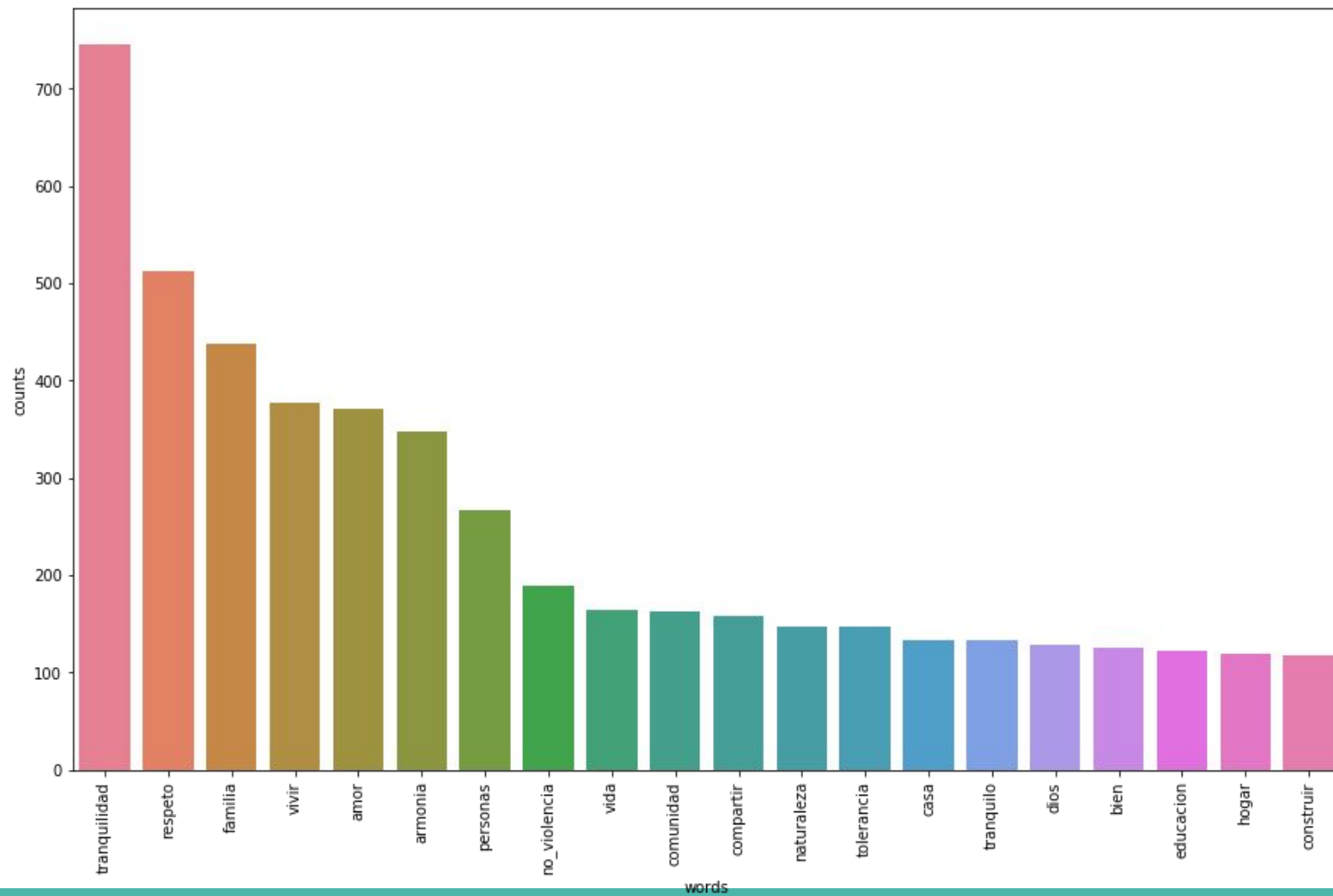
```
(0, 393)      1  
(0, 168)      1  
(0, 2)        1  
(0, 360)      1  
(0, 150)      2  
(0, 67)       1  
(0, 100)      1  
(0, 426)      2  
(1, 183)      1
```



20 most common words



20 most common words



# Encontrar mejor LDA

## Encontrar el mejor LDA

Realizaremos una búsqueda de parámetros para encontrar el mejor LDA. Para esto se tendrá en cuenta la variación de 2 parámetros:

- `n_components`, el cual es el número de topicos (grupos) que queremos formar.
- `learning_decay`, el cual controla la razón de aprendizaje.

## Grid Search

Este proceso de búsqueda de parámetros consume bastante tiempo, puesto que construye múltiples modelos de LDA para todas las posibles combinaciones de dichos parámetros en el diccionario **param\_grid** (cross-validated grid-search).

```
# Se definen los parámetros de búsqueda
search_params = {'n_components': [10, 15, 20, 25], 'learning_decay': [.5, .7, .9]}

# Se inicializa el modelo
lda = LatentDirichletAllocation()

# Se inicializa la clase GridSearchCV
model = GridSearchCV(lda, param_grid=search_params)

# Se realiza Grid Search
model.fit(data_vectorized)
```

# Escogencia del mejor modelo

```
# Best Model
lda_model = model.best_estimator_

# Model Parameters
print("Mejores parámetros del modelo: ", model.best_params_)

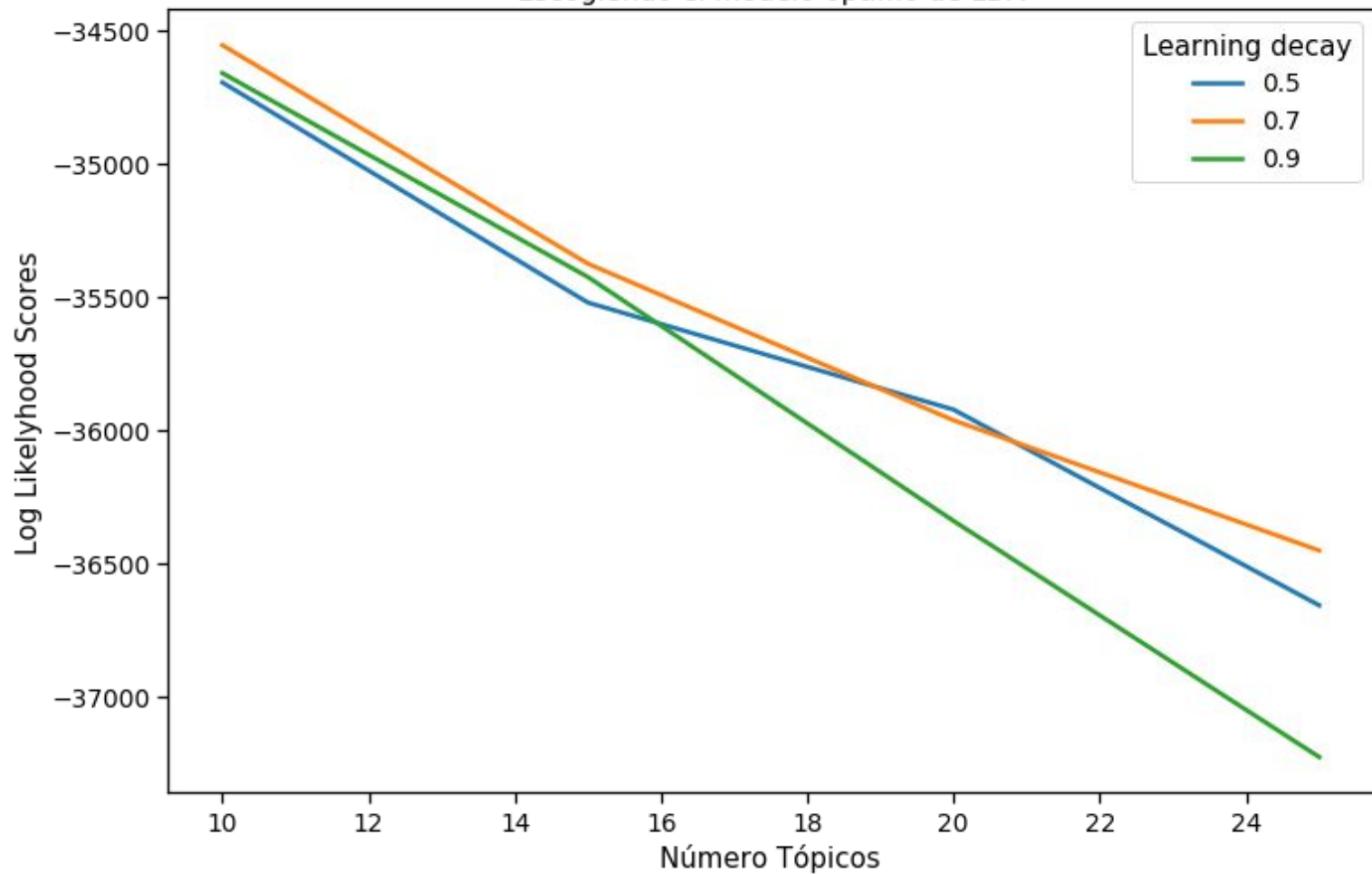
# Log Likelihood Score
print("Mejor Log Likelihood Score: ", model.best_score_)

# Perplexity
print("Perplexity del modelo: ", lda_model.perplexity(data_vectorized))
```

---

Mejores parámetros del modelo: {'learning\_decay': 0.7, 'n\_components': 10, 'random\_state': 100}  
Mejor Log Likelihood Score: -34554.234643328025  
Perplexity del modelo: 358.6854330045643

Escogiendo el modelo óptimo de LDA



# Tópico dominante en cada documento

	Topic0	Topic1	Topic2	Topic3	Topic4	Topic5	Topic6	Topic7	Topic8	Topic9	dominant_topic
Doc0	0.01	0.01	0.44	0.01	0.25	0.01	0.01	0.01	0.25	0.01	2
Doc1	0.01	0.01	0.08	0.81	0.01	0.01	0.01	0.01	0.01	0.07	3
Doc2	0.01	0.01	0.01	0.01	0.23	0.01	0.6	0.01	0.01	0.1	6
Doc3	0.25	0.01	0.01	0.01	0.46	0.01	0.01	0.01	0.01	0.21	4
Doc4	0.33	0.02	0.26	0.02	0.02	0.02	0.27	0.02	0.02	0.02	0
Doc5	0.29	0.31	0.01	0.01	0.18	0.15	0.01	0.01	0.01	0.01	1
Doc6	0.01	0.01	0.61	0.01	0.01	0.01	0.29	0.01	0.01	0.01	2
Doc7	0.03	0.03	0.03	0.03	0.41	0.03	0.39	0.03	0.03	0.03	4
Doc8	0.19	0.01	0.01	0.01	0.01	0.01	0.67	0.08	0.01	0.01	6
Doc9	0.01	0.01	0.39	0.01	0.01	0.01	0.51	0.01	0.01	0.01	6



# Distribución de los tópicos en los documentos

	Topic Num	Num Documents
0	9	474
1	4	378
2	6	368
3	2	363
4	8	273
5	3	229
6	1	221
7	5	214
8	0	213
9	7	201

# 10 palabras de cada tópico

Topic #0:

igualdad amigos entorno cuidar espacios bienestar acciones seguro construccion calle

Topic #1:

respeto personas demas\_personas diferencias contar respeto\_hacia tratar raza gente ambiente

Topic #2:

tranquilo sociedad no\_violencia vida mejor educacion oportunidades social saber siempre

Topic #3:

armonia felicidad generar hijos corazon vida sana\_convivencia alegria ser\_humano barrio

Topic #4:

tranquilidad respeto solidaridad confianza personal tolerancia reflejar valores hogar dialogo

Topic #5:

construir mundo libertad respeto armonia lograr perdonar bueno derechos territorio

Topic #6:

amor respeto tolerancia pensar conflictos personas seguridad diferente persona respeto\_tolerancia

Topic #7:

dios solo entender existir musica disfrutar escuchar projimo conflicto encuentro

Topic #8:

naturaleza casa sentir convivencia problemas representar llegar conciencia forma tener\_buena

Topic #9:

familia vivir comunidad personas armonia compartir bien interior union bueno

Selected Topic: 1

Previous Topic

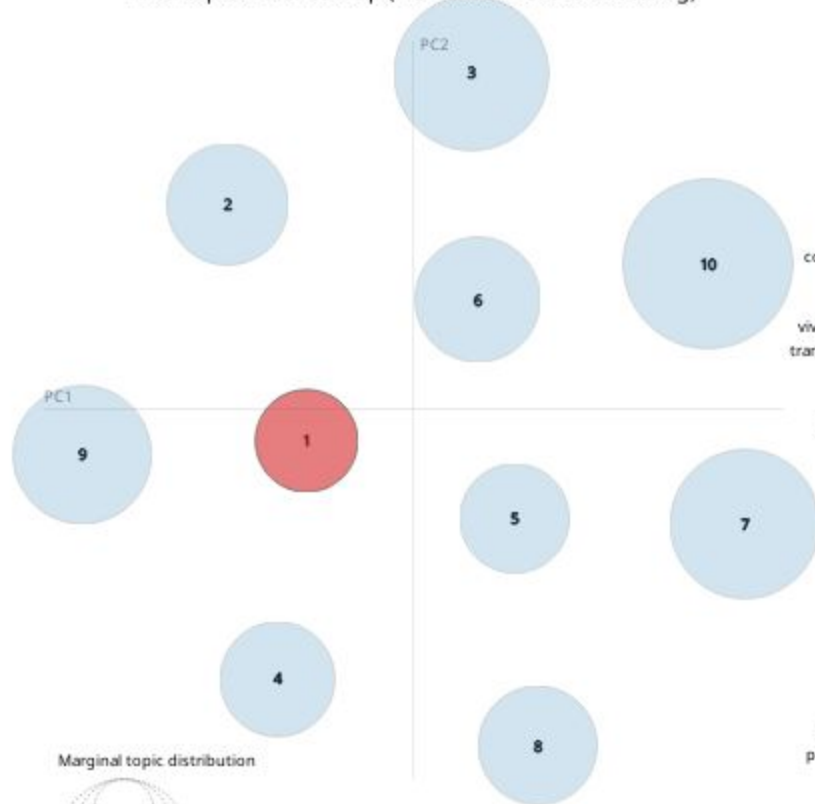
Next Topic

Clear Topic

Slide to adjust relevance  
metric:<sup>(2)</sup>  $\lambda = 1$ 

0.0 0.2 0.4 0.6 0.8 1

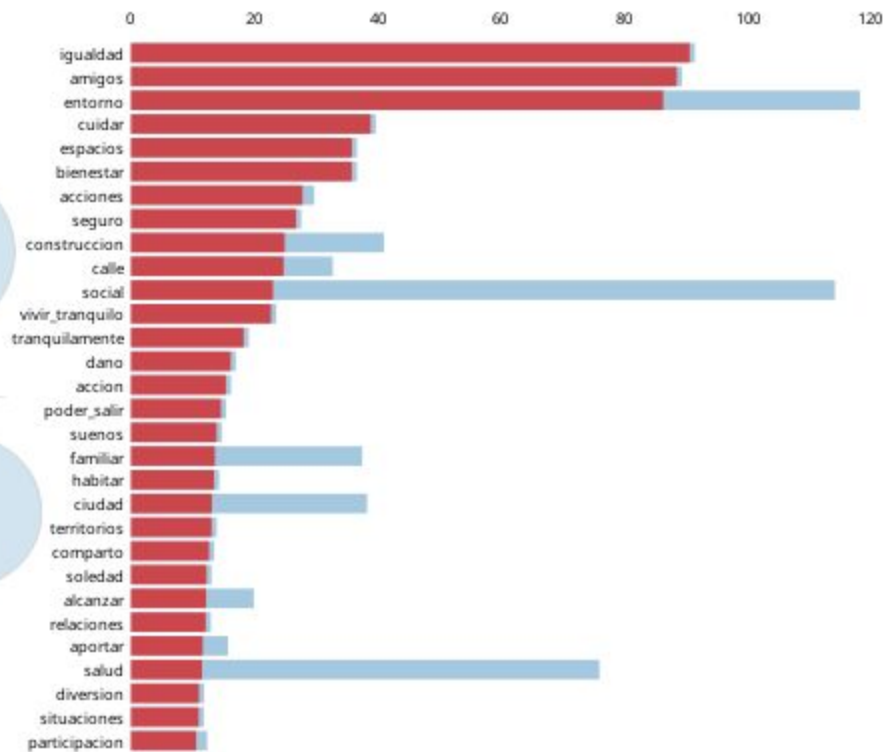
Intertopic Distance Map (via multidimensional scaling)



Marginal topic distribution



Top-30 Most Relevant Terms for Topic 1 (6% of tokens)



Overall term frequency

Estimated term frequency within the selected topic

1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * (\sum_t p(t | w) * \log(p(t | w)/p(t)))$  for topics  $t$ ; see Chuang et. al (2012)2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t)/p(w)$ ; see Sievert & Shirley (2014)

# Se exporta el modelo para usarlo en la aplicación

## Exportamos el modelo

Exportamos el modelo LDA para poder predecir futuros textos y clasificarlos en una categoría correspondiente.

```
from sklearn.externals import joblib  
joblib.dump(lda_model, 'modelo_entrenado.pkl')  
['modelo_entrenado.pkl']
```

Exportamos el modelo Vectorizer para crear nuevas matrices de dispersión acorde a los datos entrenados.

```
joblib.dump(vectorizer, 'modelo_vectorizer.pkl')  
['modelo_vectorizer.pkl']
```

## Importamos los modelos generados anteriormente

```
lda_model=joblib.load('modelo_entrenado.pkl')
```

```
vectorizer = joblib.load('modelo_vectorizer.pkl')
```

# Procesamiento de un nuevo texto

Evaluaremos un nuevo texto, pero antes tenemos que hacer una respectiva limpieza.

```
text= "Para mi la paz es tener tranquilidad en mi hogar con mi familia, mis amigas y mi perro"
```

```
text= re.sub('[,\.\\"!\\](?0-9]', '', text).lower().strip()
```

```
def sent_to_words(sentences):
    for sentence in sentences:
        yield(gensim.utils.simple_preprocess(str(sentence), deacc= True)) # deacc=True removes punctuations

data = [text]
data_tokenized = list(sent_to_words(data))
```

```
stop_words_exceptions=["no"]
def remove_stopwords(texts):
    return [[word for word in simple_preprocess(str(doc)) if (word not in stop_words) or (word in stop_words_exceptions)]

data_tokenized_nostops= remove_stopwords(data_tokenized)
```

```
# Build the bigram and trigram models
bigram = gensim.models.Phrases(data_tokenized_nostops, min count=7, threshold=5) # higher threshold fewer phrases.
```



```
4]: text_processed=vectorizer.transform(data_list)
print(text_processed)
```

```
(0, 181)    1
(0, 214)    1
(0, 426)    1
```

```
5]: topic_probability_scores = lda_model.transform(text_processed)
```

```
6]: print(topic_probability_scores)
```

```
[[0.0250002  0.025      0.02500212 0.02500561 0.51681747 0.025
  0.02500464 0.025      0.02500153 0.28316842]]
```

```
7]: np.argmax(topic_probability_scores)
```

```
7]: 4
```

```
8]: topic = df_topic_keywords.iloc[np.argmax(topic_probability_scores), :].values.tolist()
```

```
9]: print(', '.join(topic))
```

```
tranquilidad, respeto, solidaridad, confianza, personal, tolerancia, reflejar, valores, hogar, dialogo, construccion, uni
dad, medio_ambiente, adultos, generar, comun, interpersonal, practicar, normas, seres_humanos, bien, educacion, permitir,
caminar, comunitario, equidad, diferencia, diversidad, economica, derechos_sociales
```