



# الإحصاء والبيانات اللغوية

<https://www.youtube.com/watch?v=2rt2nTyFeDg&t=5s>

د. سلطان المجيلول  
@Arabic\_CL

# الإحصاء والبيانات اللغوية (لغة آر)

أولاً: ما المقصود بالبيانات اللغوية؟

ثانياً: تعلم الآلة للبيانات اللغوية (اختيار النموذج الأفضل للبيانات اللغوية) [<https://machinelearningmastery.com/>]

ثالثاً: الإحصاء (ما المشكلات؟ ما الافتراضات؟ ما المشاريع التي يُمكن أن تعالج وتحلل باستخدام لغة البرمجة آر)

- التكرار والتشتت والشيوع
- علم الدلالة والمعاني
- النحو المعجمي: الانحدار اللوجستي
- البيانات اللغوية التاريخية: تحليل التقلب الاستعمال والقمم والقعور
- التنوع اللغوي: التجميع
- التنوع اللغوي: الارتباط

رابعاً: ملفات التدريب لإعداد مشاريع في معالجة البيانات (اللغوية)

# الآر أم البايثون؟

البايثون	الآر	
المبرمجون/ علماء البيانات/ المحللون للبيانات/ الأكاديميون والباحثون/ المطورون	المبرمجون/ علماء البيانات/ المحللون للبيانات/ الأكاديميون والباحثون/ المطورون	فئة المستخدمين
pandas (manipulate data) SciPy and NumPy (scientific computing) statsmodels (explore data, estimate models, perform statistical tests) Matplotlib and seaborn (graphics) Scikit-learn (machine learning)	dplyr, tidyr, data.table (manipulate data) stringr (manipulate strings) datasets, mlbench, e1071, etc. (explore and perform statistical tests) ggplot2 (visualization) caret (machine learning)	المكتبات
مناسب لمهندسي البيانات، ومفيد لما هو أكثر من علم البيانات وتحليلها الأكواد سهلة التعلم	مناسب للإحصائيين، وأكثر متعة وفعالية في تمثيل البيانات الدالة الواحدة قد تُكتب بأكثر من طريقة (الأكواد معقدة)	حالات الاستخدام (تحليل البيانات)

# مبادئ أساسية لعمل مشاريع للبيانات اللغوية

- إتقان علم الإحصاء يمنحك مزيدًا من التمكن.
- قد يجعلك الإحصاء مثبطًا لأن الأدوات والتحليلات الإحصائية قد أصبحت أكثر تشابكًا وتعقيدًا؛ لأن التحليل الإحصائي يتضمن العديد من المسارات. نحتاج إلى اختيار بيانات لغوية مناسبة وتقنية تحليلية فعالة وتفسير مناسب للنتائج.
- لا تحتاج إلى معرفة الخوارزميات الرياضية فالحاسب هو الملاذ لاستخلاص المعلومات من البيانات اللغوية.
- لا تكن ضحية للتقنيات الإحصائية التجارية في السوق الإحصائية.
- تذكر دائما بأن الأسلوب الإحصائي الأقوى هو المنطق السليم.
- النتائج تعتمد على ما تهيئه من بيانات لغوية.
- ركز على مشروع واحد.
- المشاريع في تحليل البيانات اللغوية (العربية) عديدة. النصيحة المهمة: لا تترك مشروعًا لم تُنهه وتدخل في مشروع آخر.

Jason Brownlee

[Machine learning is taught by academics, for academics.](#)  
That's why most material is so *dry* and *math-heavy*.

[Developers need to know what works and how to use it.](#)  
We need *less math* and *more tutorials with working code*.

**Me: Let's have fun...**

**Fun makes us more optimistic in my progress**

# ماذا نعني بالبيانات اللغوية؟

- تُحول البيانات اللغوية النصية الخام المكتوبة أو المنطوقة إلى بيانات تتضمن string، و numeric، و integer، و decimal، في رأس الجدول
- تتضمن instances أو observations (صفوف) و attributive أو variables (أعمدة) في ملفات مفصولة متغيراتها وقيمها في جداول: csv، و relational database tables، و JSON، و fixed-width formatted file، و ...calc
- يتضمن ال strings حروفا أو أرقاما أو ترميزات، أما ال numeric فيتضمن أرقاما فقط (البيانات اللغوية)، أو integer.
- المتغير variable: يكون حدثا أو فكرة أو مرحلة زمنية يراد قياسها
- المتغير اللغوي: أسماء أو صفات أو ضمائر أو عبارات إلخ.

**(كلما زاد الوقت في التعلم [مستقل] زادت الدرجات في الاختبار [تابع]): الوقت يغير الدرجة والعكس غير صحيح**

1. المستقل: لا يتأثر بالتابع الخاضع للقياس: الجنس، أو العمر، أو الوعاء غير المتأثر بالتابع.
2. التابع: يتأثر؛ مثلا: التعريف والتنكير بين الذكر والأنثى، أو الصفات والأفعال والأسماء والضمائر وأدوات الربط بين الأوعية والمجالات والموضوعات.

# ماذا نعني بالبيانات اللغوية؟

## • أنواع المتغيرات

- المتغيرات اللغوية مقابل المتغيرات المستقلة (التأويلية):

1. المتغير الاسمي **nominal** / الفئوي **categorical**

2. المتغير الترتيبي **ordinal** / الرتب **ranks**

3. متغير المقياس **scale (ratio)**

- هل هناك علاقة بين جنسي المتحدثين (متغير تأويلي اسمي) واستخدام الضمائر الشخصية (متغير لغوي مقياسي)؟
- هل لكفاءة المتحدث اللغوية (متغير ترتيبي) تأثير على استخدام الكلمات الطويلة (متغير لغوي مقياسي)؟
- هل بين استخدام ضمير المتكلم واستخدام ضمير المخاطب (وكلاهما متغيران لغويان مقياسيان) علاقة؟

4. المتغير الناتج (المتغير اللغوي مثل: الخصائص النحوية) ومتغير المتنبئ (المتغير المستقل مثل: المتغيرات السياقية المرتبطة بتلك الخصائص)

# ماذا نعني بالبيانات اللغوية؟ (بنية البيانات)

The screenshot shows an Excel spreadsheet with the following data:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1		Nouns	Adjectives	Verbs	Pronouns	Coordinators									
2	AR_A01	196.322	66.10338	156.064	56.163	33.7972167									
3	AR_A02	205.915	74.47974	182.913	49.8357	27.9299014									
4	AR_A03	249.636	79.65032	192.326	31.0831	28.6546867									
5	AR_A04	197.414	79.56241	189.955	50.721	23.868722									
6	AR_A05	220.954	65.87137	168.568	51.8672	25.93361									
7	AR_A06	245.234	70.58217	185.987	41.7311	31.9422978									
8	AR_A07	221.1	77.23372	151.943	41.898	24.7349823									
9	AR_A08	197.693	75.51127	157.315	54.5359	29.8898794									
10	AR_A09	203.869	66.96429	207.341	59.5238	30.2579365									
11	AR_A10	231.361	54.83405	182.78	50.0241	19.2400192									
12	AR_A11	240.822	54.33186	173.275	30.3475	29.8580519									
13	AR_A12	199.39	61.03764	196.846	95.6256	37.1312309									
14	AR_A13	238.36	71.86235	192.308	39.9798	23.2793522									
15	AR_A14	221.253	81.62218	173.511	39.0144	24.6406571									
16	AR_A15	225.146	57.01754	193.47	78.9474	39.4736842									
17	AR_A16	234.266	59.44056	191.808	67.4326	40.959041									

Annotations in the image:

- المتغير الاسمي (Nominal variable) - Green text, pointing to column A.
- المتغير الترتيبي (Ordinal variable) - Blue text, pointing to columns B-F.
- متغير المقياس (Scale variable) - Red text, pointing to the numerical values in columns B-F.

# خوارزميات تعلم الآلة

• تعلم الآلة غير الموجه **unsupervised** والموجه **supervised**:

1. غير الموجه: التجميع.

2. الموجه: تصنيف **classification** والمخرج منه عبارة عن متغيرات فئوية، أو انحدار **regression** (الخطي وغير الخطي) والمخرج منه عبارة عن متغيرات رقمية

3. نماذج دقة التصنيف

- ✓ Linear Discriminant Analysis LDA (simple linear)
- ✓ Classification and Regression Trees (CART) <- nonlinear
- ✓ k-Nearest Neighbors (KNN) <- nonlinear
- ✓ Support Vector Machine (SVM) (nonlinear)
- ✓ Random Forest (RF) <-complex nonlinear
- Regularized Regression
- Naïve Bayes
- Decision Trees <-complex nonlinear
- Linear Regression
- Logistic Regression (mixed effects models)

## اختيار النموذج الأفضل في تعلم الآلة (التصنيف)

- عند البدء في مشروعك المحدد، قد تتساءل عن بعض النماذج، وأيها يصلح؟

- هناك أربعة أساليب لتقدير وتقييم أداء النماذج

الأول: تقسيم البيانات إلى بيانات تدريب 80% وبيانات اختبار 20%

الثاني: استخدام أسلوب التمهيدي bootstrap

الثالث: استعمال طريقة التحقق من الدقة بتقسيم البيانات إلى عدة أجزاء **k-fold cross-validation** ✓

الرابع: استخدام طريقة leave one out cross-validation [التحقق من الدقة باختبار كل instance (صف) مع بقية

الصفوف باعتبارها مجموعة بيانات تدريبية .

# اختيار النموذج الأفضل في تعلم الآلة (التصنيف)

- من الطرق التي تساعد على تقييم الخوارزمية بدوالها ومكثتها في الآر ما يعرف بـ k-fold cross-validation. يشير الـ k إلى عدد مرات تقويم البيانات، وإعطاء فكرة أولية عن دقة المتنبئات
- تقييم نتائج النموذج بالتحقق من دقة المتوسط، وذلك بتقييم كل نموذج 10 مرات [10-fold] بين بيانات التدريب والاختبار 80% وبيانات التحقق 20%

```
# Run algorithms using 10-fold cross-validation
```

```
trainControl <- trainControl(method="cv", number=10)
```

```
metric <- "Accuracy"
```

```
# LDA
```

```
set.seed(7)
```

```
fit.lda <- train(Species~., data=dataset, method="lda", metric=metric,  
trControl=trainControl)
```

```
# CART
```

```
set.seed(7)
```

```
fit.cart <- train(Species~., data=dataset, method="rpart", metric=metric,  
trControl=trainControl)
```

```
# KNN
```

```
set.seed(7)
```

```
fit.knn <- train(Species~., data=dataset, method="knn", metric=metric,  
trControl=trainControl)
```

```
# SVM
```

```
set.seed(7)
```

```
fit.svm <- train(Species~., data=dataset, method="svmRadial", metric=metric,  
trControl=trainControl)
```

```
# Random Forest
```

```
set.seed(7)
```

```
fit.rf <- train(Species~., data=dataset, method="rf", metric=metric, trControl=trainControl)
```

# اختيار النموذج الأفضل في تعلم الآلة (أر أنموذجا تطبيقيا)

• استخلاص دقة كل نموذج

Models: lda, cart, knn, svm, rf

Number of resamples: 10

## Accuracy

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
lda	0.5833333	0.6666667	0.6666667	0.6833333	0.7291667	0.8333333	0
cart	0.5833333	0.7500000	0.8333333	0.8000000	0.8958333	0.9166667	0
knn	0.5000000	0.6666667	0.7500000	0.7083333	0.7500000	0.9166667	0
svm	0.5833333	0.6666667	0.7500000	0.7416667	0.8333333	0.8333333	0
rf	0.6666667	0.7708333	0.8333333	0.8166667	0.8333333	0.9166667	0

## Kappa

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
lda	0.375	0.50000	0.500	0.5250	0.59375	0.750	0
cart	0.375	0.62500	0.750	0.7000	0.84375	0.875	0
knn	0.250	0.50000	0.625	0.5625	0.62500	0.875	0
svm	0.375	0.50000	0.625	0.6125	0.75000	0.750	0
rf	0.500	0.65625	0.750	0.7250	0.75000	0.875	0

• الغابة العشوائية rf نموذج فعال للقيم المستمرة المقيسة measured

• تحليل المحدد الخطي lda يصلح للمتغيرات المستقلة المستمرة والمتغيرات التابعة التصنيفية class أو categories (عكس ANOVA). كذلك يصلح للقيم المنفصلة أو المتقطعة أو المتميزة discrete المعدودة counted

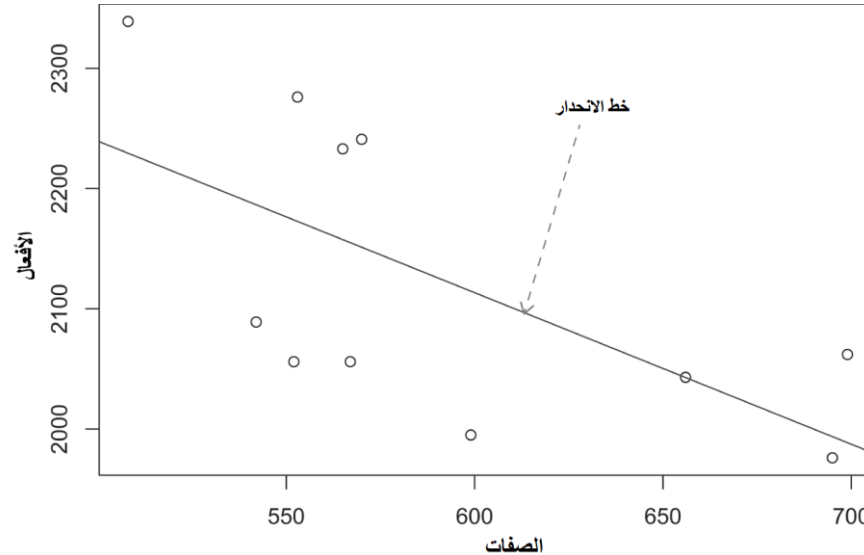
# التكرار والتشتت والتنوع

## أولاً: التكرار

- مفاهيم أساسية: التكرار المطلق، والتكرار النسبي، وأساس المعيار (في كل حجم معين)
- لتمثيل التكرار عدة أساليب بعد معالجتها. من أفضل الأساليب: خط الانحدار أو خط التوافق الأفضل

التكرار											
508	542	552	553	565	567	570	599	656	695	699	الصفة
2339	2089	2056	2276	2233	2056	2241	1995	2043	1976	2062	الفعل

- هل هناك علاقة بين استخدام الصفات والأفعال؟ الجواب: نعم، والعلاقة تناسبية عكسية (كلما زاد استخدام المؤلفين للصفات قل استخدامهم للأفعال).



# التكرار والتشتت والتنوع

- التكرار: التوزيع الطبيعي وغير الطبيعي
- القيمة الشاذة
- قياس النزعة المركزية
- القياسات الإحصائية: المتوسط، والوسيط، والمدى، والانحراف المعياري... [عمليات أولية بسيطة] والنمذجات كالانحدار الخطي والانحدار اللوجستي (نمذجة التأثيرات الممتزجة **mixed-effects models**)
- الاختبار الإحصائي: إحصائيات تتجاوز العينة للاستدلال على شيء يتعلق بالمجموع اللغوي في العينة: (مثلا: كاي تربيع وحجم التأثير)
- حجم التأثير (مثال لحالة):  
قيمة الارتباط الأسماء وعدد الصفات (  $r = .52$ ، مجال الثقة 95%  $[.46, .58]$  ) وهي ليست كبيرة كقيمة الارتباط بين الأفعال والضمائر (  $r = .81$ ، مجال الثقة 95%  $[.775, .836]$  ) التي تفسر تواردهما في ثلثي الحالات

# التكرار والتشتت والتنوع

- التشتت. مثال لتبسيط فكرة التشتت، جمع بيانات لغوية وهيئتها لأجل معرفة مدى تشتت استعمال (السعادة) في التغريدات حسب عدد المستخدمين بين المناطق الإدارية في المملكة العربية السعودية.
- التشتت يخرجنا من مشكلة Whelk: (الاعتماد على التكرار وحدة يُعد مضرًا دائمًا، وتشير هذه المشكلة إلى التوزيع غير المتكافئ). ما الحل المقترح؟ لنفترض الآتي:

الجزء 6	الجزء 5	الجزء 4	الجزء 3	الجزء 2	الجزء 1	كلمة السعادة في التواصل الاجتماعي
200,000	200,000	200,000	200,000	100,000	100,000	مجموع العينة الخام (الكلمات)
10	24	0	2	4	10	تكرار الكلمة
5	12	0	1	2	10	التكرار النسبي للكلمة في كل 100 ألف كلمة
نعم	نعم	لا	نعم	نعم	نعم	هل تشتمل على الكلمة؟

- باستخدام توزيع Juillard's D: نستخلص المتوسط (5) ثم الانحراف المعياري (4.55)، ويُقسم الانحراف المعياري على المتوسط: (الناتج: 0.91). يلي ذلك استخراج الحد الأقصى لقيمة التنوع حسب عدد أجزاء المنطقة (الجذر التربيعي لعدد أجزاء المنطقة - 1)

$$\text{توزيع جولاند} = 1 - \frac{0.91}{\sqrt{6-1}} = 0.59$$

- **الصفير = التوزيع متباين للغاية و (1) = التنوع متساوٍ للغاية**

# التكرار والتشتت والتنوع

استخدام التكرارات النسبية لكلمة لتفسير حقيقة أحجام العينات غير المتساوية

➤ التكرار النسبية (لكل مليون كلمة) لكلمة السعادة  
الفترة  
التكرار النسبي

2010	2009	2008	2007	2006
1283.21	1505.46	1623.78	1609.75	1473.69

➤ الفروق بين التكرارات النسبية  
الفترة  
القيمة1- القيمة2

2009/2010	2008/2009	2007/2008	2006/2007
222.25-	118.32-	14.03	136.06

➤ التكرار النسبي ذو اللوغاريتم المحول (لكل مليون)  
الفترة  
log2

2010	2009	2008	2007	2006
10.33	10.56	10.67	10.65	10.53

# علم الدلالة والمعاني

- قياس التوارد العشوائي (التكرارات الملحوظة والمتوقعة): تقييم لقوة معاني الكلمات من عدم القوة
- القياسات في علم الدلالة والمعاني للكلمات تتطلب معالجة للتصاحب اللفظي
  - التكرار النسبي
  - المعلومات المتبادلة
  - الأرجحية اللوغارثمية
  - قيمة-ز
  - قيمة-ت
  - الدايس
  - اللوج دايس
  - نسبة اللوج
  - الفعالية ذات النهاية الصغرى
  - دلتا ب
  - كوهن-د

# النحو المعجمي (الانحدار اللوجستي)

## بناء النموذج ثم اختيار الأكثر دقة

اسم النموذج	الناتج	المتنبات المتضمنة	النتيجة بناءً على ناتج البرنامج الإحصائي
المنطلق	النوع	[اسم]	نموذج المنطلق مع الحصر فقط
✓ الأول	النوع	نوع السياق	ذو دلالة إحصائية، أي: أفضل بكثير من نموذج المنطلق معيار معلومات أكايكي = 30.86
الثاني	النوع	نوع السياق نوع الاسم	أخطاء المعيار الكبيرة -> حدث خطأ ما يُسبب نوع الاسم مشكلة مع الفصل الكامل.
الثالث	النوع	نوع السياق طول المركب.الاسمي	ذو دلالة إحصائية، أي: أفضل بكثير من نموذج المنطلق؛ ومع ذلك، ليس أفضل بكثير من النموذج الأول معيار معلومات أكايكي = 31.91 (أكبر من معيار معلومات أكايكي للنموذج الأول)

التقدير (الأرجحية اللوغارتمية)	خطأ المعيار
-21.056	4530.376
23.889	4530.376
18.733	6706.381
18.733	9244.108
39.961	8958.692

التقدير للارجحية اللوغارتمية	خطأ المعيار	قيمة ز (اختبار زاي والد)	قيمة الاحتمال	التقدير (الأرجحية)	95% مجال الثقة أعلى	95% مجال الثقة أقل
-3.219	1.020	-3.156	0.002	0.04	0.189	0.002
6.802	1.247	5.457	0.000	900	21421.229	116.878
0.037	0.039	0.939	0.348	1.037	1.138	0.966

نوع السياق	نوع الاسم	طول المركب.الاسمي	النوع (التعريف والتكبير)
أ_غير محدد	د_عَلَم	29	ب_تعريف
ب_محدد	ب_معدود_جمع	15	ب_تعريف
ب_محدد	ج_غير.معدود	7	ب_تعريف
ب_محدد	ب_معدود_جمع	7	ب_تعريف
ب_محدد	ج_غير.معدود	27	ب_تعريف
ب_محدد	أ_معدود_مفرد	14	ب_تعريف
أ_غير محدد	أ_معدود_مفرد	4	أ_تكبير
ب_محدد	د_عَلَم	17	ب_تعريف
أ_غير محدد	أ_معدود_مفرد	51	أ_تكبير
ب_محدد	أ_معدود_مفرد	26	ب_تعريف
ب_محدد	أ_معدود_مفرد	3	ب_تعريف
ب_محدد	أ_معدود_مفرد	4	ب_تعريف
ب_محدد	أ_معدود_مفرد	6	ب_تعريف

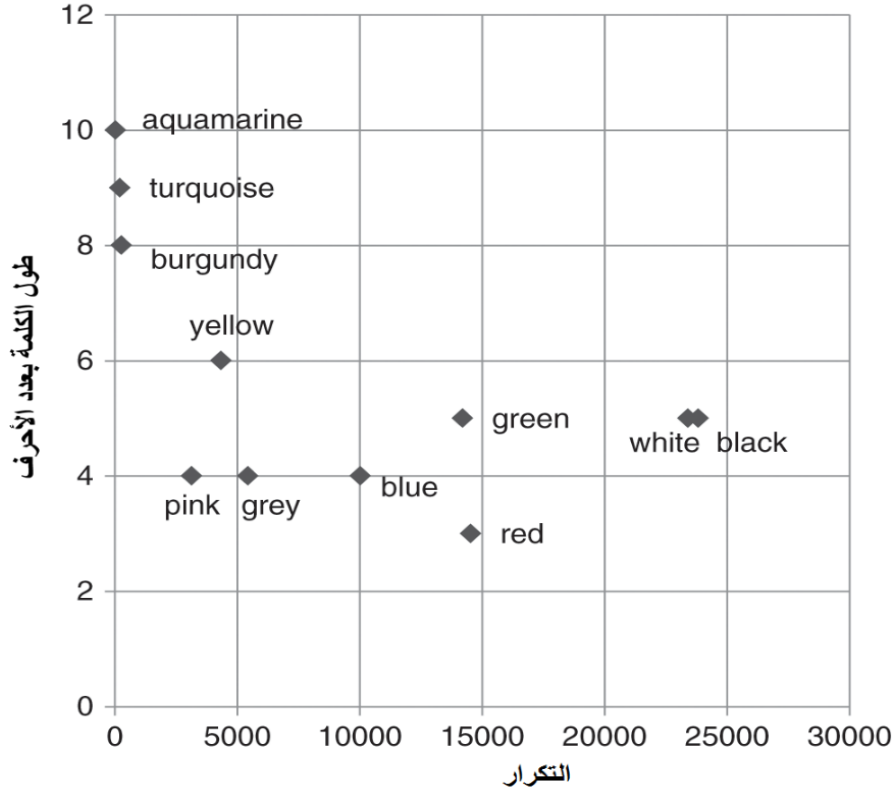
< التمييز بين المنطلق وبين القيم الأخرى للمتغيرات الفئوية (الاسمية والترتيبية): أ\_ (المنطلق)، وبقية الترميز للمتغيرات الفئوية

< المتنبات: نوع السياق، ونوع الاسم، وطول المركب الاسمي

< المتغير الناتج: التعريف والتكبير

< المتغير الناتج (المتغير اللغوي مثل: الخصائص النحوية) ومتغير المتنبئ (المتغير المستقل مثل: المتغيرات السياقية المرتبطة بتلك الخصائص)

# التنوع اللغوي: التجميع clustering



1. إحدائية التكرار وطول الجملة ل green و blue [14205, 5] و [10035, 4] على التوالي.

2. يتطلب منا اختصار المسافات بتحويل القيم إلى  $z$ -scores<sub>2</sub>

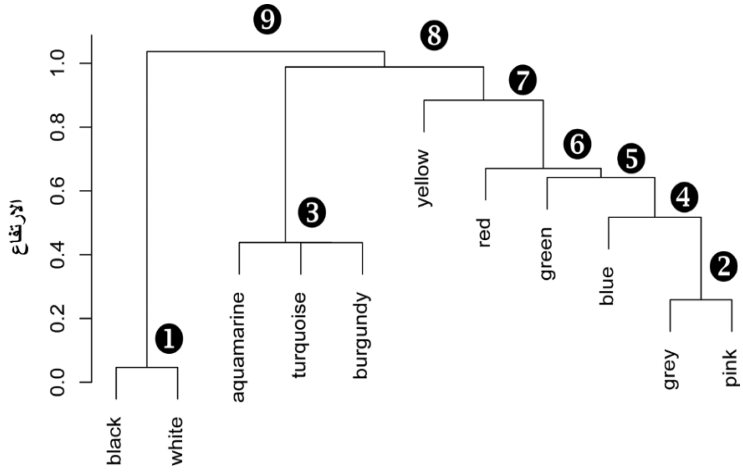
$$3. \text{قيمة-}z_2 = \frac{\text{القيمة-المتوسط}}{\text{الانحراف المعياري(العينة)}} \text{ (تختلف عن قيمة } z \text{ (z-score))}$$

4. هناك: تحتاج إلى استخراج المتوسط، ثم استخراج الانحراف المعياري [الجذر التربيعي لمجموع مربعات انحراف القيم عن المتوسط مقسومة على عدد القيم-1]

$$\text{الانحراف المعياري للعينة} = \sqrt{\frac{\text{مجموع مربعات المسافات من المتوسط}}{\text{مجموع عدد أجزاء المدونة-1}}}$$

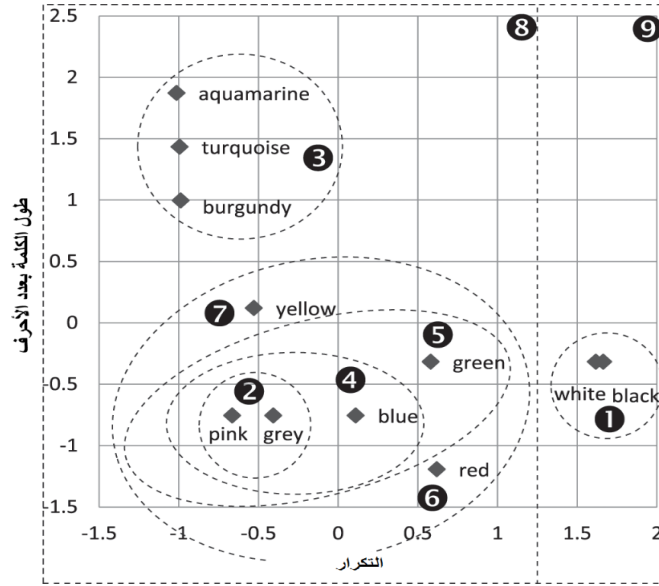
# التنوع اللغوي: التجميع clustering

إحداثيات التكرار وطول الجملة لـ green و blue [14205, 5] و [10035, 4] على التوالي  
 إحداثيات قيمة-ز2 الجديدة لـ green و blue [0.58, -0.32] و [0.11, -0.76] على التوالي

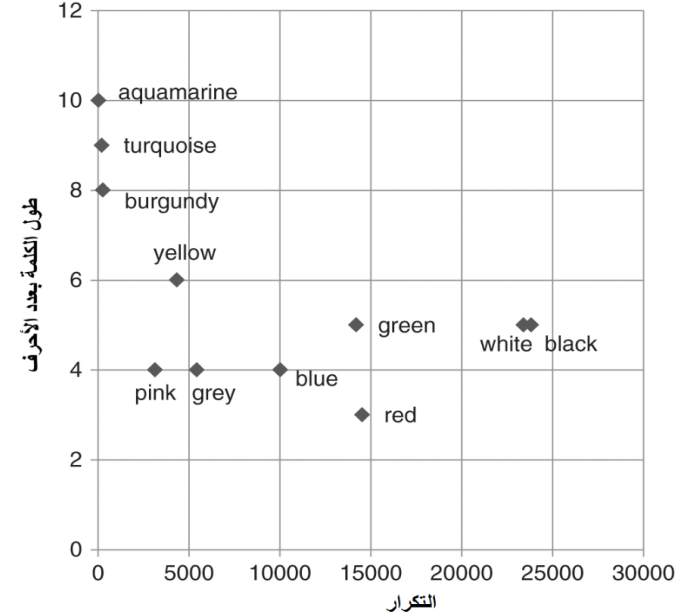


تجميع الارتفاع للمفرد

مخطط تفرعي لطريقة الربط المنفرد



التجميع



إحداثيات التكرار وطول الكلمة

$$\sqrt{(0.11 - 0.58)^2 + [-0.32 - (-0.76)]^2} = 0.64$$

المسافة الإقليدية:

$$|0.11 - 0.58| + |-0.32 - (-0.76)| = 0.91$$

مسافة مانهاتن:

# البيانات اللغوية التاريخية: تحليل التقلب الاستعمال

## تقلب الاستعمال للكلمة

تمثيل نتائج تحليل التقلب الاستعمالي UFA على البيانات اللغوية التاريخية (كلمة الحرب مثلاً)

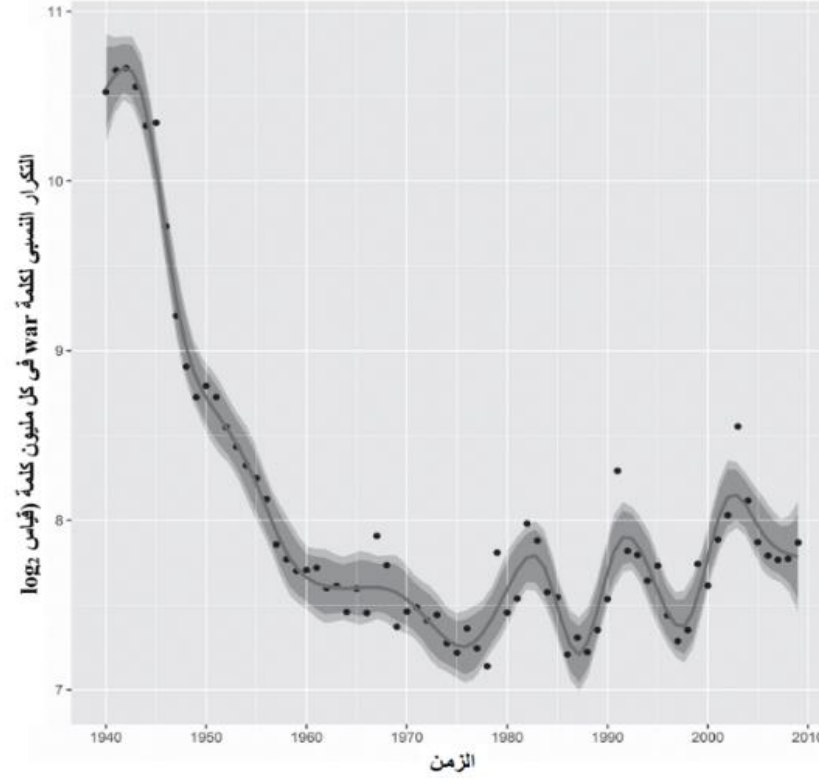
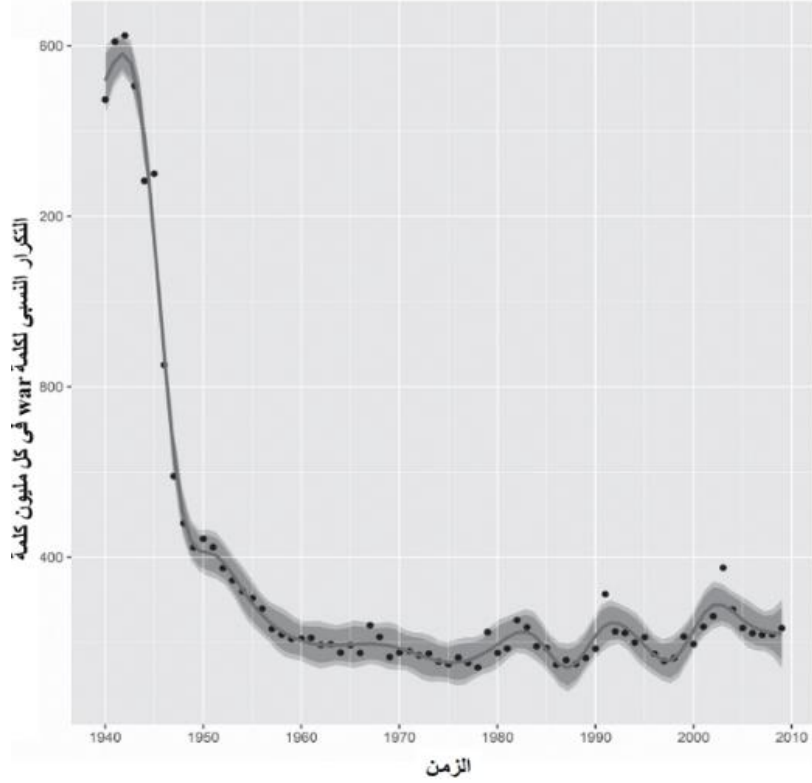
---

## القمم والقعور للكلمة والمتصاحبات

الهدف: تصنيف المتصاحبات على أنها متسقة consistent، أو مستجدة initiating، أو مهجورة terminating، أو متحولة transient

الإجراء: تمثيل نقاط الاختلاف (نقاط ذات قيم معاملات المواءمة الذي [نسبة عدد الموافقة على عدد الحالات كنقاط اختلاف للمعنى للمتصاحبات])  
بآلية القمم والقعور peaks and troughs

# البيانات اللغوية التاريخية: تحليل التقلب الاستعمال



تحليل التقلب الاستعمالي UFA على  
البيانات اللغوية التاريخية

• تمثيل لحالة (الحرب)

X = السنوات

y = التكرار النسبي

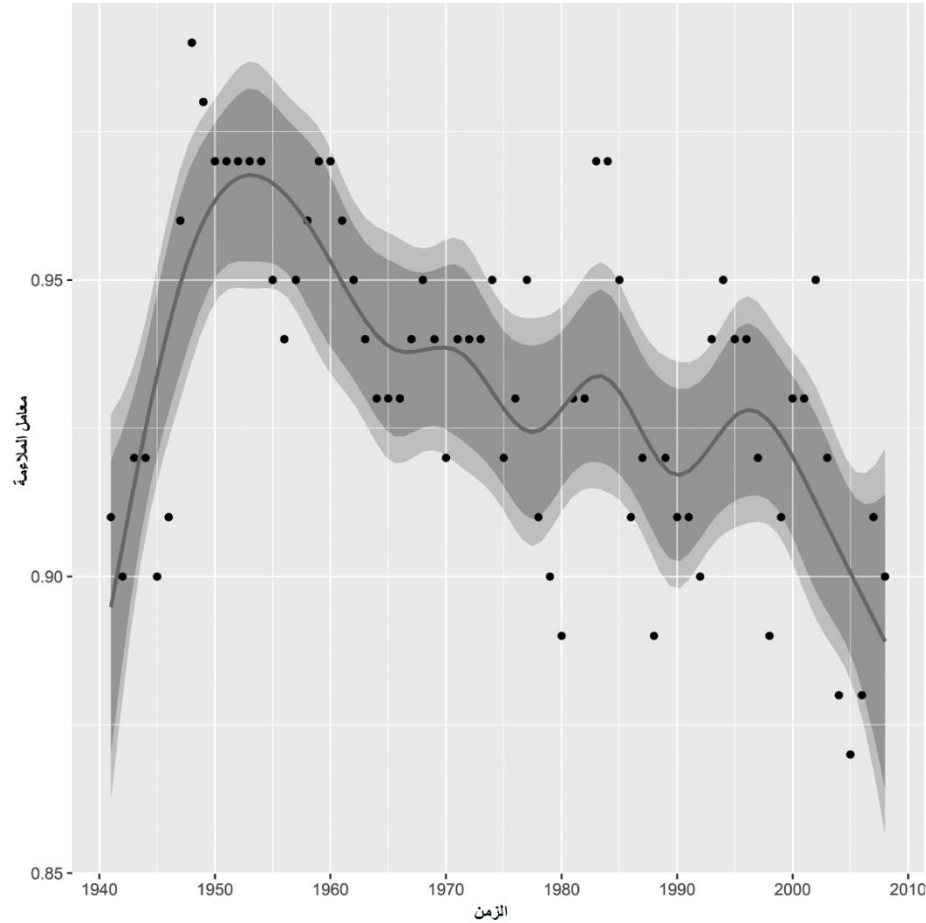
# البيانات اللغوية التاريخية: تحليل التقلب الاستعمال

## القيم والقصور للكلمة والمتصاحبات

• تمثيل لحالة (الحرب ومتصاحباتها)

$X =$  السنوات

$y =$  معامل الملاءمة أو المواءمة

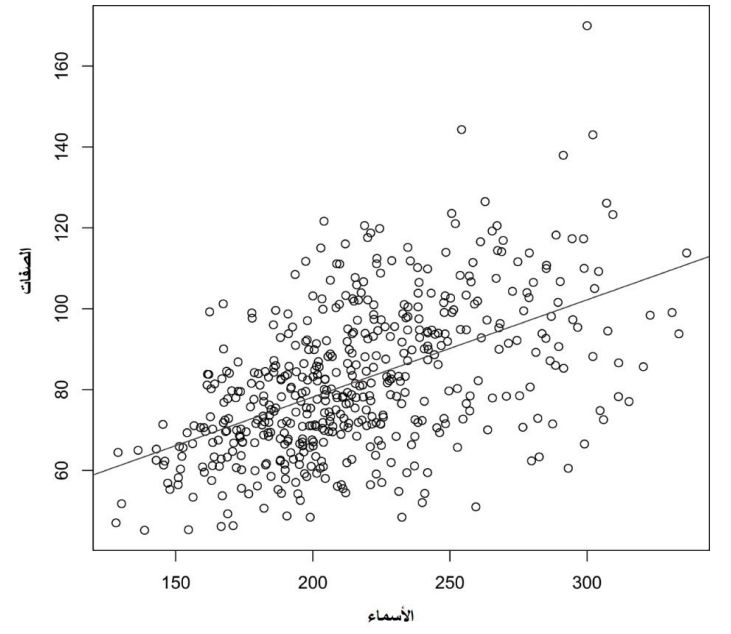
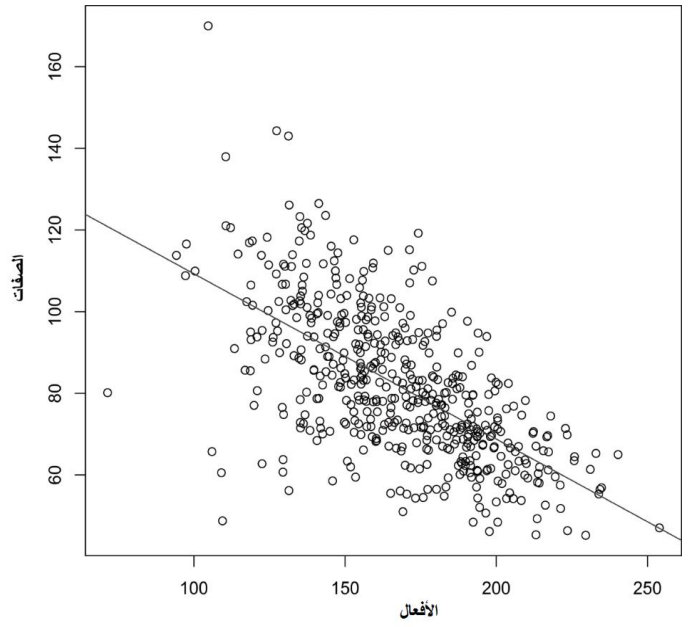
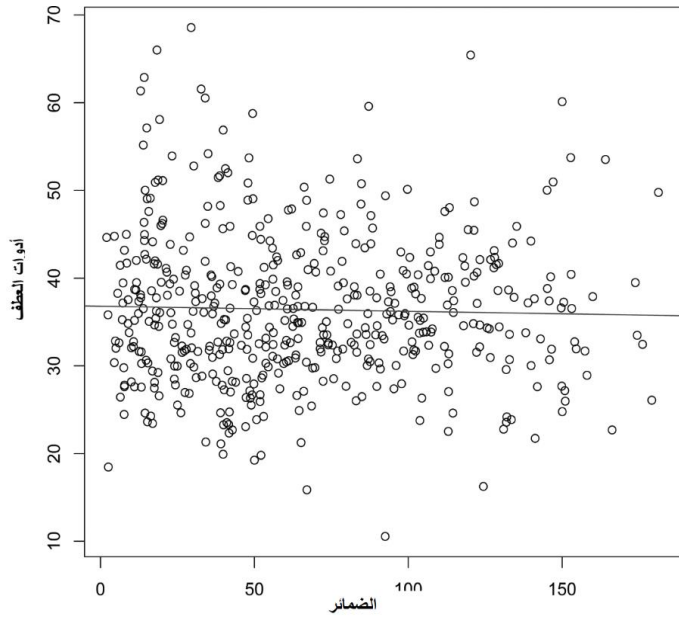


# التنوع اللغوي: الارتباط

- هل بين متغيرين لغويين ارتباط؟
- إذا كان لدينا عشر ملاحظات فقط من نصوص مكتوبة أو حوارات فنحن بحاجة إلى ارتباط كبير قدره (0.63) للوصول إلى الدلالة الإحصائية
- في حالة وجود 100 ملاحظة يكون المزيد من التغير مقبولاً وتنخفض القيمة الحرجة الضرورية لارتباط بيرسون إلى (0.2)
- في حالة وجود 1000 ملاحظة سيكون هناك ارتباط ضئيل بقيمة (0.06) وذو دلالة إحصائية
- من المنظور اللغوي فإن للارتباط بقيمة 0.06 (أو 0.2 هنا) تأثيراً عملياً ضئيلاً للغاية من جهة قوة العلاقة بين المتغيرات
- أيا كانت القيمة؛ لا بد من الإبلاغ عنها

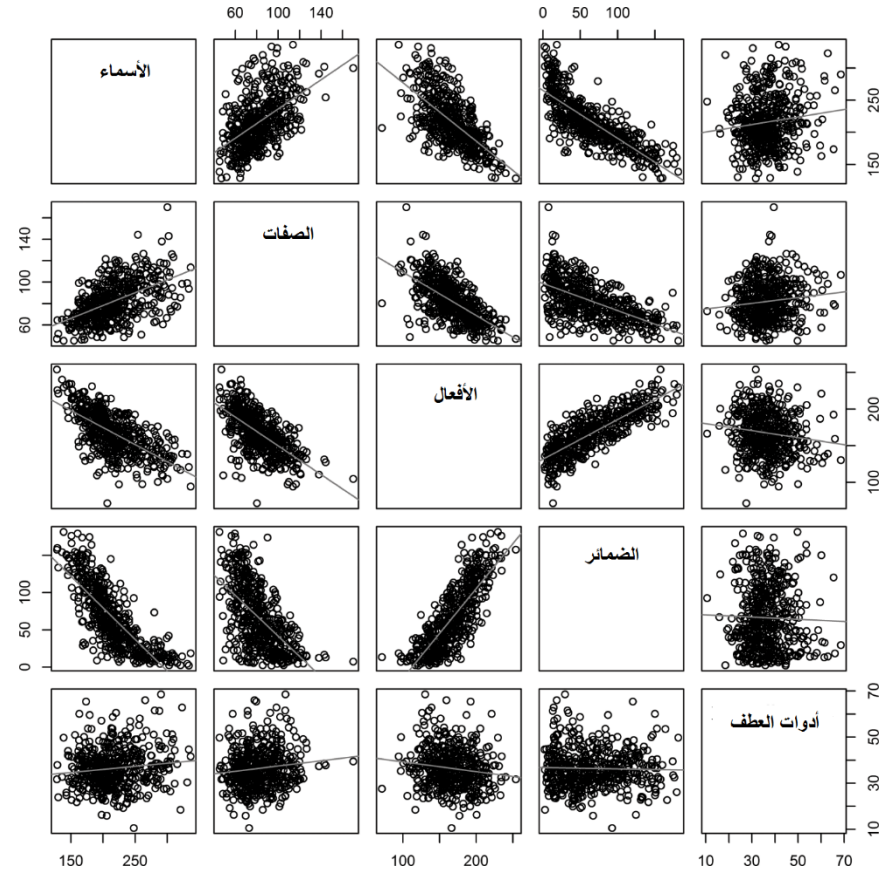
# التنوع اللغوي: الارتباط

• هل بين متغيرين لغويين ارتباط؟



# التنوع اللغوي: الارتباط

• الارتباط: مخطط التبعثر بين الأسماء والصفات والأفعال والضمائر والأدوات



# ملفات التدريب

Name	Date modified	Type
1-AnalyzeData	4/2/2020 8:45 AM	File folder
2-PrepareData	4/2/2020 8:45 AM	File folder
3-Algorithms	4/2/2020 8:45 AM	File folder
4-EvaluateAlgorithms	4/2/2020 8:45 AM	File folder
5-ImproveResults	4/2/2020 8:45 AM	File folder
6-FinalizeModel	4/2/2020 8:45 AM	File folder
7-Other	4/2/2020 8:45 AM	File folder
8-CaseStudies	4/2/2020 8:45 AM	File folder
README	2/14/2018 12:06 PM	Text Document

- التعرف على آر والمكتبات (تحليل البيانات، وتهيئة البيانات، والخوارزميات وتقويمها، وتحسين النتائج، واختيار النموذج، وتحسين أداء النموذج، ودراسة الحالات [البدء في مشروع من البداية إلى النهاية] المجلد [9].

<https://machinelearningmastery.com/>

Name	Date modified	Type
1.data_intro_tabdelimited	4/2/2020 8:45 AM	File folder
2.data_vocabulary	4/2/2020 8:45 AM	File folder
3.data_semantics&discourse	4/2/2020 8:45 AM	File folder
4.data_lexico_grammar	4/2/2020 8:45 AM	File folder
5.data_register_variation	4/2/2020 8:45 AM	File folder
6.data_sociolinguistics	4/2/2020 8:45 AM	File folder
7.data_change_over_time	4/7/2020 7:46 AM	File folder
8.data_everything_together	4/2/2020 8:45 AM	File folder
exercise_answers	4/2/2020 8:45 AM	File folder
R	4/13/2020 5:44 PM	File folder

- الإحصاءات الوصفية والاستدلالية وتحليل التكرار والشيوخ والخطاب والنحو المعجمي والتنوع واللغويات الاجتماعية والنفسية والتاريخية [8-1]

نهاية العرض

شكرًا لاستماعكم

[salmujaiwel@hotmail.com](mailto:salmujaiwel@hotmail.com): للتواصل

---