

مقدمة في معالجة اللغة العربية

سري السباعي Serry Sibae

Intern researcher at PSU in NLP

باحث متدرب في معالجة اللغة العربية في جامعة الامير سلطان

github: <https://github.com/serrysibae>

linkedin: <https://www.linkedin.com/in/serry-sibae/>

مقصدنا السؤال والنقاش ... لا كثرة الكلام والرقاش

فاسأل وناقش دونما تحرج ... فاعلم نصفه السؤال تدرج

لَكِنْ بشرطِ الدَّقِّ والتَّأدُّبِ ... فسيئُ الاخلاقِ يُحَرِّمُ يُسَلِّبِ

ومن يجد من خطأ أو سهو ... يُهْدِي جَمِيلاً دونما ترو

فلتحمزمو العقول والاذهانا ... فالرحلة خليط استبانا

الرقاش هو زخرفة اللوح وهو الحية لذلك في جلدها
[بعض الكلمات مكتوبةً عروضياً فزيد فيها الالف مثل "الاذهانا"]

فهرست المحاضرة:

١. تعريف بالمجال وأصوله وتعريفها

٢. أركان العلم الأساسية

٣. أهم أقسام العلم بالنسبة للعربية

٤. أحدث التقنيات

٥. تطبيق عملي على ثلاثة تطبيقات

٦. مفاتيح للبدء في المجال

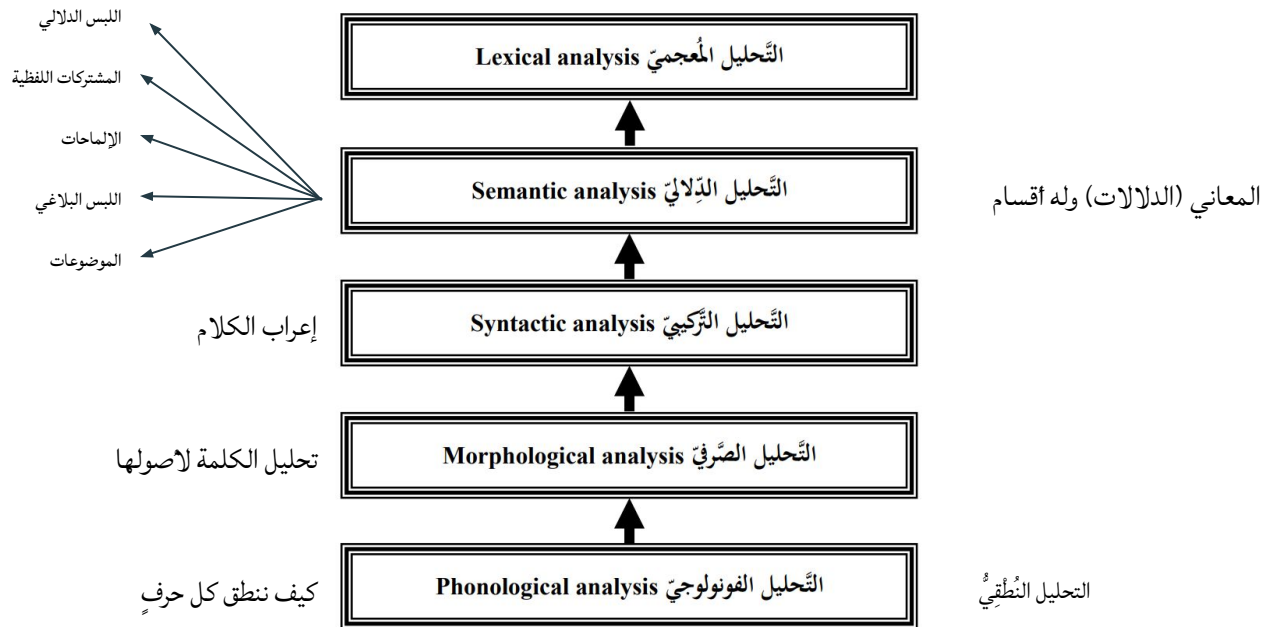
كَيْفَ نَسْتَوْعِبُ مَا نَقْرَأُ؟

الْحَمْدُ لِلَّهِ الْمُعِيدِ الْمُبْدِي ... حَمْدًا كَثِيرًا وَهُوَ أَهْلُ الْحَمْدِ

وَمَا مِثْلُهُ فِي النَّاسِ إِلَّا مُمَلَّكًَا ... أَبُو أُمِّهِ حَيُّ أَبُوهُ يُقَارِبُهُ

أي ليس في الناس حي يقاربه إلا مُمَلَّكًَا أبو أمه أبوهُ أي هشام

الحمد لله المُعِيدِ المُبْدِي ... حمداً كثيراً وهو أهلُ الحمدِ



تعريف المجال

هو مجال يَبيِّنُ (أي أنه خليطُ عدة مجالات أخرى) هدفه معالجة النصوص صوتيَّةً أو مكتوبة لإنتاج خدماتٍ مفيدة

فماذا تتوقعون هو مكوَّنٌ من ؟

1. علوم الحاسب
2. علوم اللغويات
3. الرياضيات (علمُ التَّعاليم)
4. علم تعلّم الآلة (وهو فرع من السوابق ولكن أُفردَ لبَسْطَته)

المقصد من الذكر : المقصد الاساسي من ذكر هذه العلوم هو الإلماح لعلاقة كلٍّ منها بهذا المجال لا ذكرَ تفاصيلها فهذا لا يكون إلا في محاضرات كثيرة

خلاصتها في دقائق:

برزخ نوع البرمجة والمشهور الشيئية وركناها

● الاصول سبعة:

- سير البرنامج (التصور البرمجي)
- المتغيرات
- العمليات وأنواعها
- الإدخال والإخراج
- الشروط المنطقية
- الحلقات (المكرّر)
- تتبع الأخطاء ومعالجتها

○ القولية

○ الوراثة

● الفرع التطبيقي (هياكل البيانات والخوارزميات)

○ هياكل البيانات وأصولها

■ المصفوفات وفروعها كالارتال و الاكداس

■ العقديّات (الرؤوسية) منها المتصلات والبيانات

■ المُرّمّزات مثل القاموس في مُبِشَن

○ الخوارزميات وأصولها ثلاثة

■ حساب التعقيد (الكرّبة)

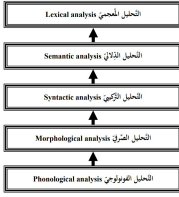
■ الترتيب

■ البحث

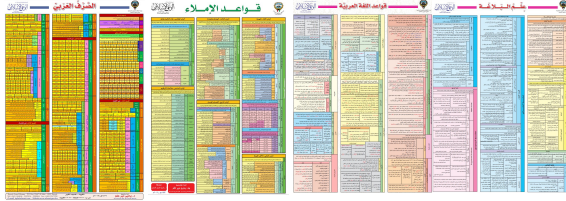
علوم الحاسب مع المجال

فن الحوسبة وعلاقته بالمجال أمان أساسيان:

1. التفكير المنطقي (البرمجي): كيف نستطيع أن نمثل أفكارنا على ما يفهمه الحاسوب وخاصة في مسألة المصفوفات والأرقام
2. الخوارزميات وهياكل البيانات: كيف يمكن أن نخزن ونسترجع البيانات ونتعامل معها بشكل سريع فعال فما زالت المعلومات تزيد أُسِّيًّا
 - a. قواعد البيانات (المتجهية) وغيرها
 - b. وحدة معالجة المرصوصات (الموترات) Tensor Processing Unit
 - c. أي لغة نختار وكيف نبرمج
 - d. تصميم خوارزميات جديدة أو تحسين الموجود



اللغويات (كلها في ورقة .)



صَرَفٌ بَيَانٌ مَعَانِي النَّحْوِ قَافِيَةٌ شِعْرٌ عَرُوضٌ اسْتِثْقَاقُ الْخَطِّ إِنِّشَاءٌ
مُحَاضَرَاتٌ وَثَانِي عَشْرَهَا لُغَةٌ تِلْكَ الْعُلُومُ لَهَا الْآدَابُ أَسْمَاءُ

- الفصحى والعامية: الاصل أن نستعمل الفصحى ولكن هذا المجال غايته التجارة والربح ولذلك استعملت العامية بأخرة وهذا ضار على المدى البعيد
- العلم هذا بشكل مختصر قائم على النظرية الوصفية (ومنها جاءت التعبيرات الإحصائية) وهناك رأي آخر لتشومسكي في نظريته التحويلية وقد نفعت في اللغات الاصطلاحية (مثل لغات البرمجة والمنطق) ولكنها لم تُجَدِ كثيرا في النماذج وتمثيلها

يرجع لسلسلة اللغة في ورقة وفيها النحو والصرف والبلاغة والإملاء في ورقة ولمن أحب فهناك ألفية لسان العرب في علوم الادب للآثاري جمع فيها عشرة علوم من علوم العرب (باط)

خلاصة الرياضيات

- المنطق الرياضي (الجبر المُجرد)
- التحليل العددي (المقاربة العددية)
- الجبر الخطي
- طوبولوجيا (الفضاءات العامة)
- الهندسة التفاضلية

لمزيد من التفاصيل عن دراسة الرياضيات البحتة انظر [المقطع](#)

الرياضيات [علم التعاليم] مع المجال

خلاصة العلم وتطبيقه في هذا الفن هو التحليل الرياضي والإحتمالات ولكن هذين مضمّنان التالي (انظر الملخصات):

- الجبر الخطي (الجخط [نحت]): وهو معالجة البيانات مصفوفياً تسريعاً وتوفيراً للوقت [كل شيء يجب أن يعرض كمصفوفة]
- الإحصاء والإحتمالات وخلاصته التوزيع الإحصائي على أي توزيع يَكُن.
- التفاضل: حيثُ يُدرّس فيه التغير وتأثيره على الدوال وتغييرها

ولمن أحب تذوق شيء من ذلك لمن بلغ مبلغاً فليُنظر ورقة ستانفورد أو بيركلي

أصول تعلم الآلة (التعالى) إلى التعلم العميق

- أنواعه الأساسية ثلاثة
 - تعلم إشرافى (س، ص)
 - توقع
 - تصنيف
 - تعلم غير إشرافى (س، س، س، س، ...، س)
 - تجميع
 - تعلم تعزيز: بيئة مع جزاء للفعل حسنا وسوءا

ولعلي أخص لكم المجال وجميع تطبيقاته في توضيح بسيط وثلاث أسس: تمثيل المدخلات و حسن العرض للسرعة (المصفوفات والمرصوصات) واستعمال النموذج المناسب

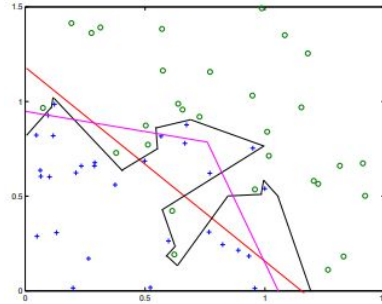
لمن أحب مقدمة قصيرة عن تعلم الآلة كتاب تعلم الآلة في 100 صفحة لطيف جدا "machine learning in 100 pages"

تعلم الآلة (التعالي) مع المجال

- ضحل: وهو قواعد فيها نوع (بعض) الذكاء وأهمها الانحدار الخطي واللوجستي وهو أساس وقاعدة التعلم العميق واستعملوا كثيرا في التصنيف خاصة
- وعميق باستعمال الشبكات العصبية وهي أشكال وألوان ولكن الذي منها للمعالجة كالتالي
 - فهم الكلمات: word2vec , BERT
 - فهم السياق
 - تصنيف الكلام
 - الترجمة (نموذج من سياق (تتابع) لآخر)
 - توليد النصوص والكلام
 - تحويل الصور لنصوص

ولعلي أخص لكم المجال وجميع تطبيقاته في توضيح بسيط وثلاث أسس: تمثيل المدخلات و حُسْن العرض للسرعة (المصفوفات) واستعمال النموذج المناسب

نظرية التعلم الإحصائي (اختياري) ؟



هنا تسقط كل التوقعات لاننا نثبت أن كل نموذج وخوارزمية لا بدّ تفشل
ولمن أحب ورقةً فليُنظر [الرابط](#)

نظرية التعلم الإحصائي هي الوحش ومن أبوابه:

- التنبؤ inference
- الصوغ formalizaion
- التعميم generalizaion
- تقارب منظم Uniform convergence
- الحساب التجريبي Empirical error
- الحجم البعدي
- التنظّم
- لا شيء بلا مقابل (نظرية لا غداء مجانيًا)

تاريخ العلم ١

لخصت تاريخ العلم من عدة مصادر ولعلي أذكر هنا خلاصته أو أكتب رؤوس أقلام مع رابط للورقة إن تيسر.

ما زال الناس يريدون الحواتل (الروبوتات) لتساعدهم وبرز منهم (يان تشي) و (بديع الزمان الجزري) و (هنري لوس وأبوه) ثم نهضت الرياضيات وملكتنا أدوات أنجعت علوم الطبيعة وظهر فيها ميكنة الحل بعد التداول والمعالجة.

ثم جاءت مرحلة اتقان الرياضيات مع غوتليب وراسل ثم كسر الأمر ونقضه غودل وظهرت نظرية الصعوبة (التعقيد) complexity وهي لب لباب علوم الحاسب.

ثم في بداية الأربعينيات وبرز نجم جون نيومان وهذا في الحرب العالمية الثانية وما بعدها. واقترب حلم الحاتل (شبيه البشر) ومن أضغاثه ظهرت صعوبة وتعقيد اللغة.

ورجا الناس محاكاة العقل البشري وظهر تورينج وتشيرش فالاول شرط الذكاء ونمذجه والثاني يسره لتبين وتدرس فعاليته بعرض الدالة من منظور حوسبي

وعام ١٩٥١ صمم أحد الباحثين ذكاء ليلعب الضامة وظهر أيضا البحث الشجري وحوسبة المنطق وهذا بين [١٩٥٧-١٩٧٤] وهو العصر الذهبي وفيه ظهرت معالجة اللغات متعمقة دلاليا مع الشبكات الدلالية وجاء مفهوم الدراسة المثالية ثم بعد الاستنتاج توضع في بيئة حقيقية (واقعية)

ومع هذا النجاح تصاعدت آمال الناس وظنوا أنهم ختموها ونسوا أنها دنيا [كمالها في نقصها] وتناقص التمويل تتابعا حتى انقطع عام ١٩٨٠ (وهنا ملاحظة أن في عام ١٩٥٨ صمم روزنبلات الشبكة العصبونية وهو مفتاح التعلم الحاسوبي)

وفي الثمانينات طوّر هينتون وروميل هارت خوارزمية تعلم الشبكات متعددة الطبقات ووقتها قامت بريطانيا واليابان بتكليف مشاريع باهظة لمعالجة اللغة

ملخص من مقال محمد عطية أحمد [حوالي ٧٠ صفحة]

انظر كتاب العربية والذكاء الصناعي

تاريخ العلم ٢ أشهر الخوارزميات تاريخيا [المصدر]

١٩٨٦	الانتشار الخلفي وخوارزمية 3iD	١٩٠١	تحليل المركبة الأساسية PCA
١٩٩٣	خوارزمية 4.5iD	١٩٢٠	المنطق الثلاثي (بدايات الضباب)
١٩٩٥	المتجهات المعتمدة وخدع الانوية وخوارزمية التجريب bs	١٩٣٧	وضوح فكرة عدم التأكد والضبابية
١٩٩٥	الغابات العشوائية وأيضاً الشبكة الالتفافية cnn	١٩٤٣	المنطق العتبي
١٩٩٧	الذاكرة العصبونية (التداخلية) RNN LSTM	١٩٥٠	اختبار تورنج
٢٠٠١	الغابات العشوائية ذات القرار	١٩٥٨	الخلية العصبية المحوسبة (روزنبلات)
٢٠٠٦	التعلم السريع	١٩٦٥	المنطق الضبابي (لطفني زاده)
٢٠٠٧	تعلم الطبقات الاناني (قطعة قطعة)	١٩٦٩	تطور الخلية المحوسبة
٢٠١٢	تدريب image net	١٩٧٢	أول تطبيق للمنطق الضبابي
٢٠١٢- إلى الآن	وانتشر المجال وبرز نجمه وتنوعت فروعها بعد ذلك	١٩٨٢	الخرائط ذاتية الترتيب

بدايات الحوسبة اللغوية العربية

١٩٧١-١٩٧٣	بدايتها مع ابراهيم انس وعلي حلمي دراسات احصائية على الصحاح واللسان والتاج وكذا عدد من الإحصائية القرآنية
١٩٨٣	مؤتمر الرباط وظهر منه كتاب اللسانية العربية
١٩٨٥	أدلة الكتابة والنظام الصرفي
١٩٨٧	محمد مرياتي نظام الاشتقاق الحوسبي
١٩٨٩	بناء معجم تركيبى وكذا نظام للتوليد والتحليل الصرفي والنحوي
١٩٩٢	الحوسبة وخصائص العربية مع بناء مُشكَّلات ومترجمات

ملخص محاضرات في اللسانيات التطبيقية [رابط](#) وهناك كتاب من أربع أجزاء لم أطلع عليه إلا سريعا

ملخص مجال معالجة اللغة

- النصوص جمعها واستيرادها وتنظيفها
- العرض (التمثيل) كيف نعرض البيانات (النصوص)
 - الكلمة والجملّة
 - العرض العدّي: في مصفوفة صفريّة أو عد نسبي أي إحصائي (خجم المصفوفة كبيرة عادةً هنا)
 - متعلم (صغير الحجم ولكنه معصور الفائدة وأشهر نموذج الكلمة المُتجهة word2vec
 - تقسيم النصوص
 - حرفا حرفا
 - كلمة كلمة
 - أجزاء
- النموذج (كيف نجعل الحاسب يتعلم العلاقات)
 - إحصائي: مثل الرمزي المنطقي (لوجستك) و حدس بايز و المتجهات المعتمدة وسلاسل ماركوف والتقسيم النوني (الكلنوني) احصائيا وهو يستخدم مع النماذج السابقة
 - نموذج عميق: الشبكة العصبونية بأنواعها تلافيفية أو ارتدادية أو بدائية أو تنبئية وغير ذلك
 - التطبيقات: وأصلها أمران توليدي أو فهمي [تَوْفُّعي و استيعابي] حيث الاول ينتج (يكتب) نصوصا والثاني يستوعب النصوص وعلاقاتها
 - ربطها باللغويات: (وهو آخر فصل من كتاب جرفاسكي) وفيه محاولة تفسير النماذج العميقة مع مقارنتها بالهياكل المعنوية [انظر شريحة الدلالة] والتراكيب النحوية وغيرها

لمزيد من التفاصيل انظر [اللبط](#)

معالجة النصوص العربية

الاصول التي ذكرها الأستاذ حبش في كتابه "مقدمة في معالجة اللغة العربية": أربعة (استخلصتها) و سَأَرَكُزُ على ما كان منها عن معالجة المكتوب وهي

1. الخط وما فيه من إعجام ورسم وتنوع الحروف فيما تستخدمه اللغات الأخرى من الحروف العربية
 - a. ويدخل تحت هذا الفصل ما تعارف عليه الباحثون من تنظيف (تهيئة) النصوص العربية قبل استعمالها في النماذج
2. الأصوات والتهجي: ولم أُطَلَّ فيه ولكن يدخل فيه كتابة الأسماء الأعجمية من لغتها للغتنا وتصحيح الأخطاء
3. الصرف وهو أطولها وفيه
 - a. التحليل الصرفي
 - b. رفع الإبهام
 - c. التقطيع
 - d. التشكيل
4. السياق ويدخل فيه معاني الدلالات

ثم بعد ذلك يكون بناء النماذج

لمزيد من التفاصيل انظر كتاب الدكتور نزار حبش وترجمته للدكتورة هند خليفة "مقدمة الى معالجة اللغة العربية"

الخط وتوابعه

وبعد أن أحسن الخطوط ، أقواه في المنسوب والمخطوط
ما وضعت أصوله القويمة ، وسلمت فروعه السليمة

وبعد إن أحسن الخطوط ... أقواه في المنسوب والمخطوط
ما وضعت أصوله القويمة ... وسلمت فروعه السليمة

تنظيف (تهيئة النصوص) وهي عملية تعارف عليها الباحثون لإزالة مالا نفع فيه ومن ذلك:

الآبيات فوق من منظومة بضاعة المجدود في الخط وأصوله للسنجاري

وبعد: إن أحسن الخطوط أقواه في المنسوب والمخطوط
ما وضعت أصوله القويمة وسلمت فروعه السليمة

- إزالة التطويلات [مرحبا , مرحبا]
- والحركات وهذه فيها إشكالات [تراجع الورقة البحثية]
- والهمزات
- والمدات
- إزالة الكلمات المعتادة المتكررة بكثرة (حروف الجر والضمائر وغيرها)

مثال برمجي

```
from camel_tools.utils.normalize import normalize_alef_maksura_ar
from camel_tools.utils.normalize import normalize_alef_ar
from camel_tools.utils.normalize import normalize_teh_marbuta_ar
```

```
sentence = "هل ذهبت إلى المكتبة؟"
```

```
print(sentence)
```

```
# Normalize alef variants to 'ا'
```

```
sent_norm = normalize_alef_ar(sentence)
```

```
print(sent_norm)
```

```
# Normalize alef maksura 'ي' to yeh 'ي'
```

```
sent_norm = normalize_alef_maksura_ar(sent_norm)
```

```
print(sent_norm)
```

```
# Normalize teħ marbuta 'ة' to heh 'ه'
```

```
sent_norm = normalize_teh_marbuta_ar(sent_norm)
```

```
print(sent_norm)
```

```
from camel_tools.utils.dediac import dediac_ar
```

```
sentence = "هَلْ ذَهَبْتَ إِلَى الْمَكْتَبَةِ؟"
```

```
print(sentence)
```

```
sent_dediac = dediac_ar(sentence)
```

```
print(sent_dediac)
```

هل ذهبت إلى المكتبة؟

هل ذهبت إلى المكتبة؟

هل ذهبت إلى المكتبة؟

هل ذهبت إلى المكتبة؟

هَلْ ذَهَبْتَ إِلَى الْمَكْتَبَةِ؟

هل ذهبت إلى المكتبة؟

الأصوات والتهجي

أغلب تطبيقات هذا القسم متعلقة بالصوتيات وليس هذا موضع درسنا ولكن دخل فيه أيضا قسما يُمكن ويحسن الإشارة لهما لتعلُّقهما بموضوعنا وهما كتابة الأسماء الأعجمية وكذا تصحيح الأخطاء الكتابة (تبعاً للكتاب)

ويجمل ذكر بعض الطرق المتبعة في تصحيح الكتابة ومن ذلك خوارزميات الكشف (heuristic) وأيضا نظام التقاسيم مثل الثنائي (Bigram)

وهذا على المستوى التركيبي

مثال برمجي

```
# Correct single word

from ar_corrector.corrector import Corrector
corr = Corrector()

all_corrections = corr.spell_correct('بختب') # return 5 corrections with top frequencies
# [(('61', 'بكتب'), ('22', 'برتب'), ('21', 'بختم'), ('9', 'بختي'), ('7', 'بخت'))]
print(all_corrections)

corr.spell_correct('من') # return true

# Correct with context

from ar_corrector.corrector import Corrector
corr = Corrector()

sent = 'أكدت قواءص التمذد فى تشاد أنها تواضضل طريقها للعاصمة'
print(corr.contextual_correct(sent))
# أكدت قوات التمرد فى تشاد أنها تواصل طريقها للعاصمة

sent = 'استتهى حدث آبل المنتظو بالإعلاخ عن مموعة من المتجات'
print(corr.contextual_correct(sent))

# انتهى حدث آبل المنتظر بالإعلان عن مجموعة من المتجات
```


الصرف

باب التصريف وفروعه من أكثر الأبواب أهمية حيث تظهر فيه معاني اللغة في تقلبياتها واشتقاقاتها ولعلي أخص الأقسام الأربعة التي ذكرناها في أول هذا القسم وهي التحليل الصرفي ورفع الإبهام والتقطيع والتشكيل

أما أولها وهو التحليل الصرفي حوسبياً بتقعيد وتأسيس ما عليه يُحلل الحاسبُ النصَّ ولنرى مثالا مكتوبا "انسألكموها" ولنحلله من جذره حتى مركباته

مثال برمجي

```
from camel_tools.morphology.database import MorphologyDB
from camel_tools.morphology.analyzer import Analyzer

# First, we need to load a morphological database.
# Here, we load the default database which is used for analyzing
# Modern Standard Arabic.
db = MorphologyDB.builtin_db()

analyzer = Analyzer(db)

analyses = analyzer.analyze('موظف')

for analysis in analyses:
    print(analysis, '\n')
```

```
diac: وُسَيِّئُونَهَا
lex: كُتِب-u_1
bw: وُ/CONJ+سَ/FUT_PART+يُ/IV3MP+كُتِب/IV4وُن/IVSUFF_SUBJ:MP_MOOD:I+هَ/IVSUFF_DO:3FS
gloss: and+_will+_they_(people)+write+it;them;her
pos: verb
root: ك.ت.ب
catib6: PRT+PRT+VRB+NOM
ud: CONJ+AUX+VERB+PRON
d1seg: وُ+_سَيِّئُونَهَا
d1tok: وُ+_سَيِّئُونَهَا
atbseg: وُ+_سَ+_يُ+_كُتِبُونَهَا
d3seg: وُ+_سَ+_يُ+_كُتِبُونَهَا
d2seg: وُ+_سَ+_يُ+_كُتِبُونَهَا
d2tok: وُ+_سَ+_يُ+_كُتِبُونَهَا
atbtok: وُ+_سَ+_يُ+_كُتِبُونَهَا
d3tok: وُ+_سَ+_يُ+_كُتِبُونَهَا
bwtok: وُ+_سَ+_يُ+_كُتِب+_وُن+_هَ
pos_lex_logprob: -3.648503
caphi: w_a_s_a_y_a_k_t_u_b_uu_n_a_h_aa
```

رفع الإبهام

بأن تعرف مكان وقسم (نوع) كل كلمة في النص هل هو فعل أو اسم وإن كان فهل هو صفة أم اسم فاعل الخ ولنرى مثالا

يا صاحب الهم إن الهم منفرج ... أبشر بخير فإن الكاشف الله

الم الم الم الم بانه ... ان ان ان ان اوانه

عِشْ اِبْقِ اسْمُ سُدْ قَدْ جُدْ مُرِ اِنَّهْ رِفِ اِسْرِ نَلْ ... غِظِ اِرْمِ صِبِ اِحْمِ اغْزِ اسْبِ رُعْ زَعْ دِلِ اِثْنِ نُلْ
وَهَذَا دُعَاءٌ لَوْ سَكَتُ كُفَيْتُهُ ... لِأَنِّي سَأَلْتُ اللَّهَ فَيْكَ وَقَدْ فَعَلَ

```
from camel_tools.tokenizers.word import simple_word_tokenize
from camel_tools.disambig.mle import MLEDisambiguator
```

```
mle = MLEDisambiguator.pretrained()
```

```
# The disambiguator expects pre-tokenized text
sentence = simple_word_tokenize('نَجَحَ بايَدَن فِي الْاِسْتِخَابَاتِ')
```

```
disambig = mle.disambiguate(sentence)
```

```
# For each disambiguated word d in disambig, d.analysis is a list of analyses
# sorted from most likely to least likely. Therefore, d.analysis[0] would
# be the most likely analysis for a given word. Below we extract different
# features from the top analysis of each disambiguated word into separate
# lists.
```

```
diacritized = [d.analysis[0].analysis['diac'] for d in disambig]
pos_tags = [d.analysis[0].analysis['pos'] for d in disambig]
lemmas = [d.analysis[0].analysis['lex'] for d in disambig]
```

```
# Print the combined feature values extracted above
for triplet in zip(diacritized, pos_tags, lemmas):
    print(triplet)
```

```
from camel_tools.tokenizers.word import simple_word_tokenize
from camel_tools.disambig.mle import MLEDisambiguator
from camel_tools.tagger.default import DefaultTagger
```

```
mle = MLEDisambiguator.pretrained()
tagger = DefaultTagger(mle, 'pos')
```

```
# The tagger expects pre-tokenized text
sentence = simple_word_tokenize('نَجَحَ بايَدَن فِي الْاِسْتِخَابَاتِ')
```

```
pos_tags = tagger.tag(sentence)
```

```
print(pos_tags)
```

مثال برمجي

→ ('نَجَحَ', 'verb', 'نَجَحَ')
 ('بايَدَن', 'noun_prop', 'بايَدَن')
 ('فِي', 'prep', 'فِي')
 ('اِسْتِخَابَاتِ', 'noun', 'اِسْتِخَابَاتِ')

→ ['verb', 'noun_prop', 'prep', 'noun']

تقسيم (تقطيع) الكلام

وهذا من أهم أقسام هذا المجال وينبغي عليه أغلب التطبيقات الحديثة وله أنواع كثيرة سنشرح منها أهم ثلاثة:

1. تقسيم مبني على الحروف
2. وآخر على الكلمات نفسها
3. وأشهرها على التجزيء فتصير الكلمات أجزاءً مقطعةً

وسنشرح أشهر خوارزمية لكل واحدة إن شاء الله

التقسيم على الحروف

وهو مبني على أن نعتبر الحرف هو المركب الأساسي للكلمة فمثالا:

كلمة "ملعب" إذا حللناها تصير [م,ل,ع,ب] وهكذا

وهذه الطريقة لها فائدة كبيرة إذ لا يعيقها شيء عن تقسيم أي معطى كيفما كان ولكن يعيها طول مدخلاتها وقلة بياناتها

مثال برمجی

```
# we build this by our selfs
```

```
chars = "ابتجخدذرزسسظعففكلمنهوي'!<×÷|{~:/,&,آإءةؤى"
```

```
stoi = { ch:i for i,ch in enumerate(chars) }
```

```
itos = { i:ch for i,ch in enumerate(chars) }
```

```
encode = lambda s: [stoi[c] for c in s] # encoder: take a string,  
output a list of integers
```

```
decode = lambda l: ''.join([itos[i] for i in l]) # decoder: take a
list of integers, output a string
```

```
print(encode("أَنْ يُطَلَّعَ عَلَى صُورَةٍ حَقِيقَةٍ"))
```

[58, 30, 25, 50, 29, 28, 30, 16, 1, 30, 23, 41, 18, 30, 29, 18, 23, 59, 29, 15, 32, 27, 11, 30, 57, 42, 29, 7, 30, 21, 28, 21, 28, 1, 30, 57, 42]

التقسيم على الكلمات

وهو أن نعتبر كل كلمة قسما مميزا كما هي بدون أي تغيير على أجزائها

كلمة "ملعب" تبقى "ملعب" و "ملعبا" تبقى نفسها

وهذه الطريقة تنفع حيث أن لكل كلمة متجها خاصا بها فيقل عدد المدخلات وتزيد دلالاتها ولكن يعيبها عيب قاتل وهو كبر حجم قاعدة الكلمات المستعملة وهذا سيؤدي إلى صعوبات في الحسابات وأيضا المقسم لن يستطيع فهم الكلمات الجديدة إن لم تكن موجودة قبل.

مثال برمجي

```
# we can build this by our self
```

```
text = ""
```

هذا النص هو مثال لنص يمكن أن يستبدل في نفس المساحة، لقد تم توليد هذا النص من مولد النص العربي، حيث يمكنك أن تولد مثل هذا النص أو العديد من النصوص الأخرى إضافة إلى زيادة عدد الحروف التي يولدها التطبيق.

إذا كنت تحتاج إلى عدد أكبر من الفقرات يتيح لك مولد النص العربي زيادة عدد الفقرات كما تريد، النص لن يبدو مقسماً ولا يحوي أخطاء لغوية، مولد النص العربي مفيد لمصممي المواقع على وجه الخصوص، حيث يحتاج العميل في كثير من الأحيان أن يطلع على صورة حقيقية لتصميم الموقع.

ومن هنا وجب على المصمم أن يضع نصوصاً مؤقتة على التصميم ليظهر للعميل الشكل كاملاً، دور مولد النص العربي أن يوفر على المصمم عناء البحث عن نص بديل لا علاقة له بالموضوع الذي يتحدث عنه التصميم فيظهر بشكل لا يليق.

هذا النص يمكن أن يتم تركيبه على أي تصميم دون مشكلة فلن يبدو وكأنه نص منسوخ، غير منظم، غير منسق، أو حتى غير مفهوم لأنه ما زال نصاً بديلاً ومؤقتاً.

```
""
```

```
un_set = set(text.split()+[" ", "<جهول>"])
```

```
word_indx = {word:index for index,word in enumerate(un_set)}
```

```
indx_word = {index:word for index,word in enumerate(un_set)}
```

```
# put the unknown part
```

```
def encoder(text):
```

```
    collected = []
```

```
    for word in text.split():
```

```
        if word in word_indx.keys():
```

```
            collected.append(word_indx[word])
```

```
        else:
```

```
            collected.append(word_indx['<جهول>'])
```

```
    return collected
```

```
decode = lambda x: " ".join([indx_word[indx] for indx in x])
```

```
encoded = encoder("من أن يطلع على صورة حقيقية لو يستطيع الطالب أن يفهم مفهوم اللعبة")
```

```
print(encoded)
```

```
print(decode(encoded))
```

[17, 70, 55, 33, 30, 45, 51, 51, 51, 70, 51, 8, 51]
من أن يطلع على صورة حقيقية <جهول> <جهول> <جهول> أن <جهول> مفهوم <جهول>

التقسيم المُجزَّء (الجزئي)

وهو أشهر أنواع المقسمات وهو الذي يُستعمل في أغلب إلا تكن كلُّ التطبيقات الحديثة من المحادثات (تشات جبت chat gpt) وغيرها من النماذج الكبيرة وهو مبني على مبدأ لطيفٍ جداً وهو أن تعرض الكلمات بأجزاء كل منها وكمثالٍ عليه فكلمة "ملعب" مثلاً قد تصير "مل" و "عب" فتخيل كم عدد الكلمات التي أجزاءً منها السوابق

وأشهر هذه المقسمات BPE "مقسم بأيّتي" وهدفه كان ضغط البيانات وهذا في الثمانينيات ولكن استخدم في معالجة اللغة بأخرة وظهر ويزغ وما زال. ولعلي أشرح الخوارزمية بسرعة

وعليه بُنيَ عدد من المقسمات الجديدة التي فيها تحسينات إحصائية مثل "جزء الجملة" Sentence Piece

صاح	ساح
صح	سم
صح	سن
صد	سو
صر	سي
صع	سيي
صع	شوي
صف	شوي
صق	شا
صك	شيب
صل	شيت
صم	شيج
صن	شيج
صو	شيج
ضا	شك
ضاب	شك
ضاح	شوي
ضاح	شوش
ضاح	شوط
ضد	شع
ضر	شعب
ضع	شوق
ضع	شك
ضف	شك
ضق	شك
ضل	شم
ضم	شمن
ضن	شه
ضو	شوي
	شوي

```
import re, collections
```

```
def get_stats(vocab):
    pairs = collections.defaultdict(int)
    for word, freq in vocab.items():
        symbols = word.split()
        for i in range(len(symbols)-1):
            pairs[symbols[i], symbols[i+1]] += freq
    return pairs

def merge_vocab(pair, v_in):
    v_out = {}
    bigram = re.escape(' '.join(pair))
    p = re.compile(r'(?!\S)' + bigram + r'(?!\S)')
    for word in v_in:
        w_out = p.sub(' '.join(pair), word)
        v_out[w_out] = v_in[word]
    return v_out

vocab = {'l o w <w>': 5, 'l o w e r <w>': 2,
         'n e w e s t <w>': 6, 'w i d e s t <w>': 3}
num_merges = 10
for i in range(num_merges):
    pairs = get_stats(vocab)
    best = max(pairs, key=pairs.get)
    vocab = merge_vocab(best, vocab)
    print(best)
```

r .	→	r.
l o	→	lo
l o w	→	low
e r .	→	er.

function BYTE-PAIR ENCODING(strings C , number of merges k) **returns** vocab V

```

V ← all unique characters in C           # initial set of tokens is characters
for i = 1 to k do                       # merge tokens til k times
     $t_L, t_R$  ← Most frequent pair of adjacent tokens in C
     $t_{NEW} \leftarrow t_L + t_R$            # make new token by concatenating
     $V \leftarrow V + t_{NEW}$                # update the vocabulary
    Replace each occurrence of  $t_L, t_R$  in C with  $t_{NEW}$  # and update the corpus
return V

```

Figure 2.13 The token learner part of the BPE algorithm for taking a corpus broken up into individual characters or bytes, and learning a vocabulary by iteratively merging tokens. Figure adapted from [Bostrom and Durrett \(2020\)](#).

نبدأ بأصغر مكونات النصوص: ا ب ت ث ... هـ و ي ١, ٢, ٣, ٤, ٥, ٦, ٧, ٨, ٩

بعد اللفة الأولى: اب , اج , ... , هو , هي

بعد اللفة الثانية:

بعد لفات معدودة

يتكون عندنا مجموعة أجزاء نستعملها و الأجزاء المذكورة أعلاه مثال من مقسم أرابيرت العربي

مثال برمجي من مكتبة أرابيرت ٢

```
ids = tokenizer("من أن يطلع على صورة حقيقية لو يستطيع الطالب أن يفهم مفهوم اللعبة")
['input_ids']
# tokenizer.tokenize("من أن يطلع على صورة حقيقية لو يستطيع الطالب أن يفهم مفهوم اللعبة")
tokenizer.convert_ids_to_tokens(ids)
cod = tokenizer.encode("مرحبا بك في لعبتنا")
tokenizer.decode(cod)
print(cod)
```

```
ar_toks = tokenizer.tokenize("السلام عليكم ورحمة الله وبركاته")
['السلام', 'علي', 'كم', 'و', 'ر', 'حم', '##', 'الله', 'و', 'ير', 'كات', '##']
```

Word Piece

لمزيد من التفاصيل يُرجع للمقالة التالية ([انط](#))

وهذا فيه إضافة حسابات إحصائية

$$\mathcal{O}_{\text{ML}}(\theta) = \sum_{i=1}^N \log P_{\theta}(Y^{*(i)} | X^{(i)}) .$$

Sentence Piece

وهذه طريقة أخرى تستهدف اللغات التي لا تفصل كلماتها بالمسافات المعروفة

Bigram

```
["b", "g", "h", "n", "p", "s", "u", "ug", "un", "hug"],
```

"hugs" could be tokenized both as ["hug", "s"], ["h", "ug", "s"] or ["h", "u", "g", "s"]. So which one to choose? Unigram saves the probability of each token in the training corpus on top of saving the vocabulary so that the probability of each possible tokenization can be computed after training. The algorithm simply picks the most likely tokenization in practice, but also offers the possibility to sample a possible tokenization according to their probabilities.

Those probabilities are defined by the loss the tokenizer is trained on. Assuming that the training data consists of the words x_1, \dots, x_N and that the set of all possible tokenizations for a word x_i is defined as $S(x_i)$, then the overall loss is defined as

$$\mathcal{L} = - \sum_{i=1}^N \log \left(\sum_{x \in S(x_i)} p(x) \right)$$

الدلالة

الدلالة: علاقة بين تعبير ومعنى وهي في العربية أصل وقرائن خمسة ملحق وصوتي وتركيبى وصرفي وسياقي

وأصنافها: لفظ وإشارة وعقد (حساب) وخط وحال

وعلم الدلالة هو علم المعنى ومحاوره ثلاثة ماهية وتعليلية وتدليلية ويحتوي على الكلمة والمفهوم [وهو مركب ذهني من قواعد و مؤثرات]

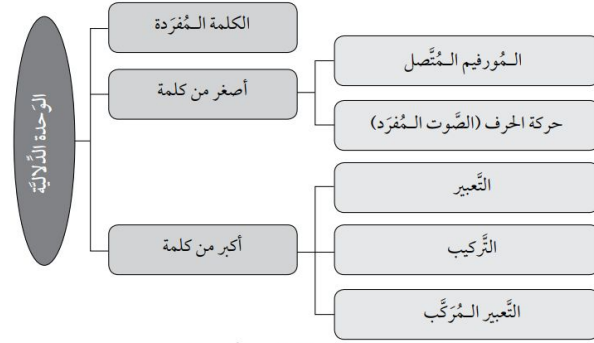
وأما الانطولوجيا: فهو مصطلح إغريقي يعني علم الوجود [وقد عرّبه للهيكلة المعنوية وغيري لـ النماذج المعرفية]

وأما هندسة [وهي قياس الواقع على الأصول أو المفاهيم] الانطولوجيا: فهي فن بناء هياكل مساعدة على ترتيب البيانات وسهولة فهمها واستخدامها واسترجاعها

الوحدة الدلالية

تلخيص لبعض المصطلحات المنتشرة:

- مورفيم: وهي وحدة المعالجة النبوية أي أصغر وحدة لها معنى
- الجذع: هو مجرد السوابق واللواحق
- والجذر: أصل الكلمة المشتقة منه
- والزوائد سوابق ولواحق [أوائل وأواخر]^(١)



الشكل ٦: أنماط الوحدة الدلالية في اللغات الطبيعية

م	أنماط الوحدات	نماذج الوحدة الدلالية
١	الكلمة المفردة (المورفيم المفرد)	إنسان
٢	أصغر من كلمة	المورفيم المتصل (س) - يعمل = التسويف
	حركة الحرف	كُتِبَ (سَ) = خطاب المفرد المذكر
٣	أكبر من كلمة	التعبير
	التعبير	فَرَسَ أَخَاسًا فِي أَسَدَاسٍ = تحوَّرَ
	التعبير المركب	تَابَعَ شُرَا = العلمية
		الحافلة النهرية

الجدول ٩: أنماط الوحدات الدلالية في اللغة العربية

(١): مصطلحيات هذا هو المشهور ولكني وجدت احسن كلمة لها لاصية: لاحقة أو إضافة حرفية تُراد في أول الكلمة أو آخرها أو تُدرج

في وسطها لتغير من دلالتها اللغوية، أو لتخصص معناها في منظومة اصطلاحية معينة، فإذا كانت في أول الكلمة تُسمى **الراعية**،

بمعنى السابقة، فالرُغف في اللغة التقدم والسبق، وإذا كانت في آخر الكلمة فهي **الكاسية**، والكُسع في اللغة الإتيان، والاصل إن

يُقال القافية إلا أن هذا المصطلح أستهلك في الشعر، أما في وسط الكلمة فهي القاجمة من القُحم وهو الدخول في الشيء.

شرح النموذج

عادة الآن صارت النماذج عصبونية شديدة العمق والتعقيد كثير عدد المتغيرات جدا ولكن هذا لا يعني أن نترك ما يُعطي نتائج جيدة وتكلفة تدريبه أقل بكثير

- نموذج إحصائي: مثل حُدس بايز و المتجهات المعتمدة و الانحدار المنطقي الرمزي وسلاسل ماركوف وغيرها
- نموذج عميق: الشبكة العصبونية بأنواعها تلاففية أو ارتدادية أو بدائية أو تنبئية وغير ذلك

[تنبيه على الاصل في الهندسة التفاضلي]

لمن أحب مراجعة مصدرٍ صعب ولكنه مفيد في تبیین أصول التعلم العميق معمما ومؤصلا كل النماذج لباب واحد فليُنظر المحاضرة و ورقتها

ولمن أحب البدء في تعلم الآلة فعليه بفهم علم البيانات أولا من numpy, pandas, matplotlib ثم البداية في الخوارزميات الرئيسية ثم ينطلق الشخص إلى التعلم العميق هناك كتاب بالعربي للاستاذ طعيمة ([رابط](#)) وله أيضا كتاب في التعلم العميق وآخر في مشاريع علوم البيانات ([رابط](#) مجموع الكتب) [ملاحظة كل كتب الدكتور مترجمة] وهناك كتب ومحاضرات كثيرة بالكلزية (eng) منشورة على الشابكة لمن أحب

كيف المجال الآن ؟

انزاح المجال تماماً نحو نماذج المحولات من ٢٠١٧ فما أن يخرج نموذجٌ حتى يتبعه آخر فظهر الآن نماذج لكل فرع من معالجة اللغة

عليكم بـ hugging face

في اكتشاف المُسمّيات الكيانية

NER ← XLM-Roberta

لتوليد النصوص

Text Gen: GPT, LLAMA, Falcon, BERT, T5 (full)

موضوع **الفعالية** في الاستخدام من أمثله التكميم (تصغير عدد الفواصل)

دخول المحولات في المجالات الأخرى كمعالجة الصور والصوتيات والحمض النووي

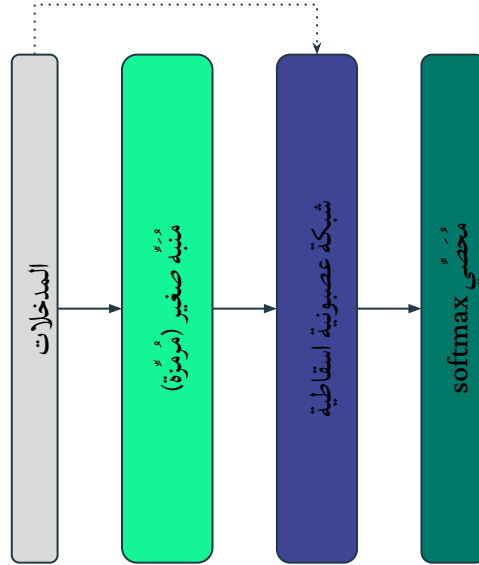
لمزيد من التفاصيل انظر كتاب Natural Language Processing with Transformers

تطبيق عملي

1. بناء شبكة انتباه صغيرة
 2. باحث معنوي
 3. استعمال متمم (chat GPT) مع هندسة المُدخلات
- إضافات إذا سمح الوقت: نموذج احصائي ساذج و انحدار منطقي رمزي (لوجستي)

لخص ما يحسن من دروس cs124 عند كتابة مسودة الكتاب

النُبيه الصغير



معادلات النموذج

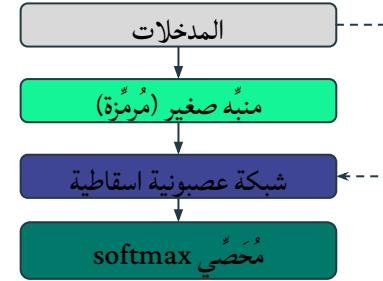
$$س_{ن \times ل} ص^{(ع)} \text{ حيث } ع = \{م, ق, ا\}$$

$$انتباه(م, ق, ا) = مَحَص(ا, م) \sqrt{\text{ابعاد}}. ق = خ_{ن \times ك}$$

$$خ = ش = ن_{ن \times م}$$

$$مَحَص(ن) = ج$$

لمزيد من التفاصيل حول حسابيات المحولات انظر [الورقة](#)



يمكننا أيضا أن نضيف طبقة ناظرية norm لمشكلة تلاشي المُشتق

حواشي توضيحية:

- سين مدخل والهزمة تُبدل بحسب اسم المصفوفة (كلها نفس الابعاد ولكن الاسم يختلف) حيث ميم هي المفتاحية وقاف هو القيمة والالف هي الاستردادية (انظر شرحها في آخر)
- فُجْداء س مع ع يعطينا ق, ا, م حيث ابعاد كل واحدة (ن × ك)
- نتيجة الانتباه مصفوفة خ_{ن × ك}
- ثم الإسقاطية مصفوفة ش_{ك × م}
- عملية التحصية على الأعمدة م في النهاية
- هذه العملية بدون نظام الموازيات المُدخلية (بدون مرصوصات)
- ج متجهة عمودي فيه احتمال الكلمة التالية وهو المُستمثل
- السهم المُنقَط إضافة ممكنة بأن تقفز للطرف التالي وقد استُعملت في الورقة

الباحث المعنوي

معادلة جاكارد يسيرة جدا ان شاء الله

ب(ج١، ج٢) = تقاطع(ج١، ج٢) \ اجتماع(ج١، ج٢)

ويمكن تحسين هذه المعادلة بعمل تقطيع نوني (كلنوني) (كلمتان أو ثلاثة الخ) وكذا يمكن تحسينها بتنظيف المدخل **تجديعا** أو **تجديرا**

أما المعادلة الأخرى فهي تكرار الكلمة و معكوس التكرار المستندي (تك متم)

تك(كلمة, مستندها) = تكرارها الى كلماته

[ت\ك]

متم(الكلمة, المستندات) = لوغ (عدد المستندات \ عدد ما جاءت فيه)

[ل(ع\ظ)]



الصورة مربوطة
فاضغطها

استعمال متمم (مدردم [نحت من مدردش ومتمم])

أصول التعامل معه أربعة:

أمور قد تفيدك في الدراسة:

- محددات إحصائية (عادة لا نستعملها) مثل وسع أو مدى العشوائية وكمية الاختيارات البحثية وغيرها
- تدريبية: كإدخال قطعة مع سؤال عنها أو سياق لفهمه وغير ذلك
- نوع المخرج نصا أو برمجة أو غير ذلك
- تأطير المخرج
- اطلب نظام ٢٠-٨٠ باريو
- اطلب خطة لدراسة موضوع أو أكثر تداخلا
- خطة للمراجعة التكرارية للضبط
- تصور المفاهيم المعقدة (طلب تصوير منطقي)
- اطلب أسئلة للتأكد من الفهم
- حول المفهوم الى قصة
- اطلب الشرح خذي خمس

هناك [bard](https://www.perplexity.ai) من جوجل وكذا موقع <https://www.perplexity.ai> ولمن أحب فهذا تلخيصي لمحاضرات الدكتور أنيس وهواري عن الموضوع [دايط](#)

خاتمة

كيف يبدأ شادي المجال فيه:

- مكتبات شديدة النفع والإفادة Useful libraries

- Camel tools
- farasa (py)
- spacy
- nltk
- gensim
- hugging face
 - transformers and tokenizers
 - datasets
- spark NLP
- Regex
- pytorch

1. تعلم البرمجة ببايثون وأتقنها ليسهل عليك برمجة النماذج
2. تعلم أساسيات الجبر الخطي والإحصاء لتفهم أصول عمل النماذج (طبعاً هنا لا حَدُّ أعلى ولكن أدنى)
3. تعلم مكتبات معالجة اللغات مثل camel-tools, nltk ومكتبات الذكاء الصناعي ويكفيك: pytorch, sklearn
4. طبق ثم طبق ثم طبق وابنِ ما استطعت من المشاريع وتجعلها تتابعاً بل تأزياً (مع بعض)
5. اقرأ إذا أحببت ما تيسر من الكتب وبالعربي كتب مركز الملك عبد الله
6. كمقدمات مع تطبيقات ومشاريع مفيدة انظر كتب الاستاذ علام طعيمة (انظر الشريحة رقم ٤٠)

لمن أحب البدء في المجال فعليه بدروس وكتاب جرفاسكي وصاحبه مارتن، [ابط](#) وهذا رابط [مقاطعها](#)

مصادر أوراق بحثية عربية

<https://github.com/iwan-rg/ArabicSurvey>

[رابط](#) لمنتقياتي من كتب مركز الملك عبد الله لمن أحب

من أهم المصادر في الترجمة "معجم مصطلحات الرياضيات" طبعة مجمع دمشق ([رابط](#))

مقالة فيها كتب لرياضيات التعلّؤ [رابط](#)

جميع الشفوات البرمجية في هذه المُسودة لمن أحب تشغيلها ودراستها ([رابط](#)) ولمن يجد خطأً فليُكرمنا به ولم يبقَ ناقصاً منها إلا برمجة النُبيه لعل الله ييسر ذلك قريباً

الشرائح من إعداد سري السباعي

Email: serrymrss@gmail.com

github: <https://github.com/serrysibae>

linkedin: <https://www.linkedin.com/in/serry-sibae/>