

Comp150-NLP Final Project

by Alena Borisenko

°_ ✧ ∖ (° ∇ °)_ ✧ _°

Idea

— — —

- Look at popular Japanese Twitter accounts
- See what words and characters are used
- Compare to popular Kanji lists and/or literature
- (Maybe) Generate a sentence
- Recommend ways to learn

Why Japanese?

Currently learning the language myself

もし日本語が読める人がいたら、教えてください！

Curious to see how Japanese complicates things

- Regular expressions? (surprisingly easy)
- Tokenizer? (unsurprisingly hard)
- Looking cool? (easier than ever)

Data

Data: Twitter Accounts

Data: Twitter Accounts



@hajimesyacho

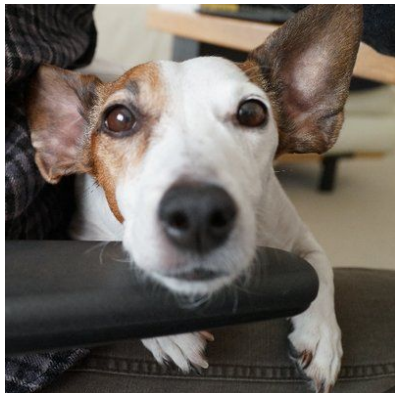
[YouTuber]
[24]

Data: Twitter Accounts



@hajimesyacho

[YouTuber]
[24]



@itoi_shigesato

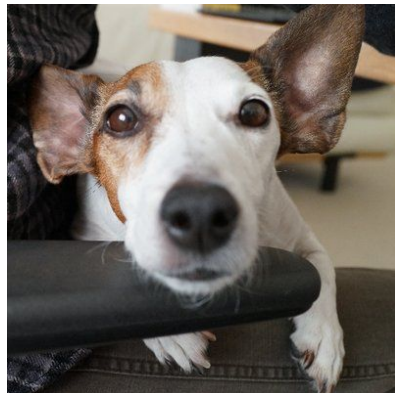
[Writer, Lyricist]
[69]

Data: Twitter Accounts



@hajimesyacho

[YouTuber]
[24]



@itoi_shigesato

[Writer, Lyricist]
[69]



@pamyurin

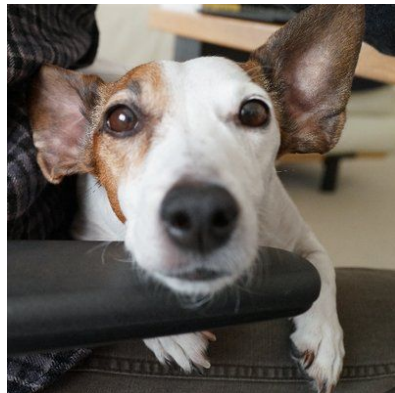
[J-Pop Singer]
[24]

Data: Twitter Accounts



@hajimesyacho

[YouTuber]
[24]



@itoi_shigesato

[Writer, Lyricist]
[69]



@pamyurin

[J-Pop Singer]
[24]



@nhk_news

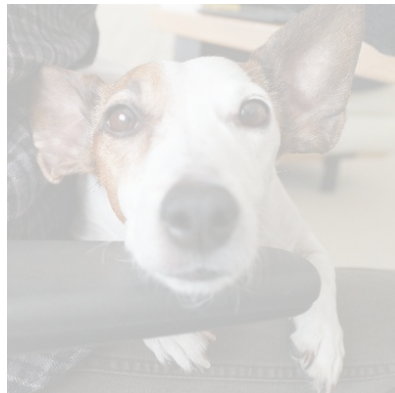
[News Org]
[N/A]

Data: Twitter Accounts



@hajimesyacho

[YouTuber]
[24]



@itoi_shigesato

[Writer, Lyricist]
[69]



@pamyurin

[J-Pop Singer]
[24]



@nhk_news

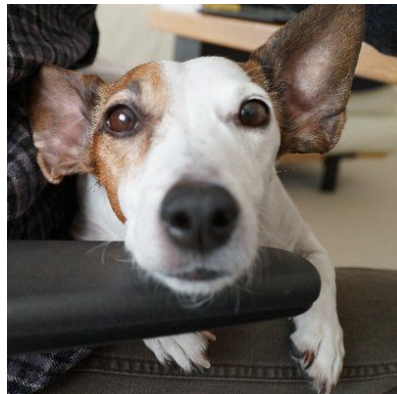
[News Org]
[N/A]

Data: Twitter Accounts



@hajimesyacho

[YouTuber]
[24]



@itoi_shigesato

[Writer, Lyricist]
[69]



@pamyurin

[J-Pop Singer]
[24]



@nhk_news

[News Org]
[N/A]

Data: Kanji

tangorin.com

Lists of common kanji

JLPT (N5 to N1)

- Japanese Language Proficiency-Test
- For non-native speakers (like TOEFL)

Jōyō (G1 to G7)

- Taught in primary and secondary school
- Issued by the Japanese Ministry of Education

Tangent: JLPT Learning Curve is Insane

JLPT N5

— — —

日 一 国 人 年 大 十 二 本 中 長 出 三 時 行 見 月 後 前 生 五 間 上 東 四 今 金 九 入 学 高 円 子 外 八 六 下 来 氣 小 七 山 話 女 北 午 百 書 先 名 川 千 水 半 男 西 電 校 語 土
木 聞 食 車 何 南 万 每 白 天 母 火 右 読 友 左 休 父 雨

JLPT N4

— — —

会同事自社発者地業方新場員立開手力問代明動京目通言理体田主題意不作用度強公持野以思家世多正安院心界教文元重近考画海壳知道
集別物使品計死特私始朝運終台広住真有口少町料工建空急止送切転研貸堂鳥飯勉冬屋茶弟牛魚兄犬妹姉漢
歌貰恵囃週室歩風紙黒花春赤青館屋色走秋夏習駅洋旅服夕借曜飲肉貸堂鳥飯勉冬屋茶弟牛魚兄犬妹姉漢

JLPT N3

— — —

政 議 民 連 対 部 合 市 内 相 定 回 選 米 実 関 決 全 表 戦 經 最 現 調 化 当 約 首 法 性 要 制 治 務 成 期 取 都 和 機 平 加 受 続 進 数 記 初 指 權 支 産 点 報 濟 活 原 共 得 解
交 資 予 向 際 勝 面 告 反 判 認 参 利 組 信 在 件 側 任 引 求 所 次 昨 論 官 增 係 感 情 投 示 変 打 直 両 式 確 果 容 必 演 歳 争 談 能 位 置 流 格 疑 過 局 放 常 状 球 職 与 供
役 構 割 費 付 由 說 難 優 夫 收 断 石 違 消 神 番 規 術 備 宅 害 配 警 育 席 訪 乘 残 想 声 登 易 助 劳 例 然 限 追 商 葉 佗 勵 形 景 落 好 退 頭 負 渡 失 差 末 守 若 種 美 命 福 望
非 観 察 段 横 深 様 財 港 識 呼 達 良 候 程 満 敗 値 突 光 路 科 積 留 他 処 努 精 太 客 散 否 師 婚 喜 浮 絶 幸 押 飛 殺 号 単 座 破 弘 庭 完 徒 勤 責 捕 危 給 苦 迎 欠 園 更 具 刻 贊 辞 因 抱 馬 愛 恐 息 遠 寒 願 髮 亡 絵
冷 適 婦 寄 込 互 束 似 幾 遊 夢 君 閉 緒 折 草 暮 酒 悲 晴 掛 到 寝 暗 盜 吸 陽 御 齒 忘 雪 吹 娘 誤 洗 慣 礼 窓 昔 貧 怒 祖 泳 杯 疲 皆 鳴 腹 煙 眠 怖 耳 頂 箱 晚 寒 忙

JLPT N2

党 協 總 區 領 県 設 改 府 查 委 軍 団 各 島 革 村 勢 減 再 稅 營 比 防 補 境 導 副 算 輸 述 線 農 州 武 象 域 額 欧 担 準 賞 迎 造 被 技 低 復 移 個 門 課 腦 極 含 蔵 量 型 況 針
專 谷 史 階 管 兵 接 細 効 丸 湾 録 省 旧 橋 岸 周 材 戸 央 券 編 搜 竹 超 並 療 採 森 競 輕 介 根 販 城 歷 將 幅 般 貿 講 林 裝 諸 劇 仏 奧 紅 詞 胃 造 航 築 雙 腰 疊 被 鉄 貨 混 純 躍 机 復 禁 昇 翌 冊 濯 移 印 池 快 勇 塔 個 逆 血 片 械 灰 門 換 温 敬 菜 沸 課 久 季 惱 珍 菓 腦 短 星 泉 皮 湖 枯 極 油 永 皮 湖 枯 含 暴 著 漁 荒 貯 虫 貝 蔵 輪 誌 荒 貯 虫 貝 量 占 庫 貯 硬 刷 符 型 植 刊 硬 刷 符 況 清 像 埋 湯 憎 針 倍 香 柱 溶 皿

JLPT N1

氏攻闢邸驅沿脚陳壞莊疎晉吳扇斗巽愁朔侃勁
統崎葬縮透妙潮憶漫諾仰穗凡腸蘭癖杜樓伽娑菖
保督避還津唱梅潛玄雷剛壯憩槽迅愉深彬畎洵旦
第授司属壁阿尽梨粘漂疾堤媛慈肖寅韵匡抄爾棕
結催康慮稻索僕仁悟懷征飢溝楊鉢礁廉眉爽耗昴紬
派及善梓仮誠桜克舖勘碎傍恭伐朽乃謹欽黎甦胤胤
案憲逮惠裂斐滑岳妊裁謡疫刈駿殼洲瞳薪情銃胤胤
策離迫露敏懇孤概熟拐嫁累睡漬享屯湧揭窰莞凜
基激惑冲是俳炎拘旭駄謙痴錯糾秦樺炊姻嵯萌碩聒
価摘崩緩排柄賠墓恩添后搬弁亮茅嫻窰縉汪有聒
提批聰需堅麻鋼須往斜菌瘳穀坪沙巖醜縉汪有聒
挙批聰需堅麻鋼須往斜菌瘳穀坪沙巖醜縉汪有聒
応郎脱射訳李頑偏豆鏡鎌桐陵紺輔擬升栓偃晏莉
企健級購芝浩鎖雰遂聡巢寸霧娛媒塀殉翠溜伎汰
検盟博揮綱剂彩遇狂浪頻郭魂椿鷄唇煩鮎允朕瑤
藤從締充典瀨摩諮岐亜琴尿弊舌禪陸姻閑凹蒔綸耶
沢修救貢賀趣勵狹卓緯詐棚吐舶峽胴劾艶鯉且榔
裁隊執鹿扱陷縦卓緯詐棚吐舶峽胴劾艶鯉且榔
証織房却願斎輝龜培壇潔鷹窮厘挿峻租蕙遥晨丞
援拡撤端弘貫蓄糧衰點酷鷹窮厘挿峻租蕙遥晨丞
施故削賃看仙軸簿艇魔宰賣掌峰圭椎詠棧隼瑛燦壺
井振密獲訟慰巡稼牧怪紫寂腐麗毛推詠棧隼瑛燦壺
護并措郡戒序稼牧怪紫寂腐麗毛推詠棧隼瑛燦壺
展就志併祉旬瞬殊淡曙辰鐘奥蓮陪侮婿表附頌昂
態異載徹營兼砲殖抽紋霞憾悅弔刮鑄裴逐卯箇枉
鮮猷陣貴欽聖噴艦披卸伏猪刁乙譜抹斐斥但楓熙
視敵我衝奏旨誇輩廷奮緇縹磁曆尼悠規嬌某詔茜戢
条維為焦勸即祥穴錦欄俗磁曆尼悠規嬌某詔茜戢
幹浜抑奪騷柳牲奇准逸漠弥宜邇淑隸某詔茜戢
独遺幕災閤舍秩慢暑涯邪昆盲衡帆船禍因阜凌濤峻
宮壘染浦甲偽帝鶴磯拓晶粗粹薰曉蝶魁雖皓勾濤
率邦奈析繩較宏謀獎眼墨訂辱狎傑酪虹惟洸晟捷
衛素傷讓鄉霸唆暖浸獄鎮羊殺羊楠莖鴻佑羅娑杓
張遣択称揺詳阻昌刺尚洞庄轄款節逆於黨緋杓
監抗秀納免抵泰拍胆彫履傘履閔鈴逆於黨緋杓
環模徵樹既霄賄朗織穩劣敦弦偵奴汽起渥鯛丙
審雄彈挑薦茂撲寬駒顯那騎稔喝錠珥臻漸懂伶颯
義益償誘隣犧堀覆虛巧毆寧窒敢拳匿蚊宵邑茄
訴緊功紛華旗菊胞靈予娠循炊胎翔襟窠妄倣勺
株標拋至範距絞泣帳垣泰忍洪醉遷蚩厄惇碧恕
姿宣秘宗隱雅綠隔悔欺憂急摂憤拙蕉藻脩啄蒔
閤昭拒促德飾唯淨諭釣朴如飽豚侍寡祿甫穰瑚
衆庖刑慎哲網膨没慘萩亭寮冗遮尺琉孟酌酉遵
評伊塚控杉奄矢暇虐肅淳祐桃扉峠痢嫡蚤悽瞭
影江致智积詩耐肺翻栗怪鵬狩疏雉庸堯嬉儉虞
松僚繰握己緊塾貞墜愚鳩鉛朱赦肇朋嚇蒼柚虞
撃吉尾宙妥翼漏靖沼嘉醇珠渦窃渴坑已暉爾柵
佐盛描俊豪敵猛飼肥架穫苗枢瑞雌暢只詢謁
核皇鈴錢豪敵猛飼肥架穫苗枢瑞雌暢只詢謁
整臨盤洪熊魅芳陰徐鬼佳獸碑又亨揜頤肢采斤
融踏項銃滯嫌懲銘糖庶潤哀鍛慨堪畔霜檀紗嵩
製壞喪操微斉剣随搭稚倬跳刀紡叙遼硝凱賦捺
票債伴携隆敷彰烈盾滋乏匠鼓恨酢頰勒替眸蓉
涉興養診症擁棋丁稿淹煮赴蛇猶扶通橘杏梓式製
響源懸託暫園丁稿淹煮赴蛇猶扶通橘杏梓式製
推儀街撮忠酸恒丹軌姫桑澄塊戲嶺杏梓式製
請創契誕倉滅揚啓依誓桂澄塊戲嶺杏梓式製
器障掲侵彦罰冒也妨把髓僧弓忌喬擗凰嗣諄治
士繼躍括肝礎之丘擦鯨呈盆亘膜奔漆苑慧絢笙壘
討筋棄謝喚腐倫棟鯨呈盆亘膜奔漆苑慧絢笙壘

Back-end

Tweepy

— — —

Simplifies use of Twitter API

Handles authentication

Search limitations:

- 200 tweets per search
- 7 day window
- ~3200 recent tweets for a single user
- Limited number of requests

Solution: get_tweets.py

Set up tweepy

Get recent tweets from api

Write an .xml file

tweets.xml

— — —

```
<tweets>
  <tweet>
    <num>1</num>
    <id_str>938251171617128449</id_str>
    <user>pamyurin</user>
    <lang>ja</lang>
    <text>
      鍵が壊れてトイレに閉じ込められて大変だった、、ふう( ͡° ͜ʖ ͡° )
    </text>
  </tweet>

  <tweet>
    <num>2</num>
    <id_str>937951688958263296</id_str>
    <user>pamyurin</user>
    <lang>ja</lang>
    <text>
      ○💎 https://t.co/hihYrB1m7L
    </text>
  </tweet>
</tweets>
```

get_kanji.py

Read in the .csv kanji guides from tangorin

Get the kanji

Write to separate file for later use

project.py

Set up Flask

Use a regular expression to isolate kanji

Use a tokenizer to isolate words (WIP)

Compare results to JLPT/Jōyō kanji lists

Generate a sentence (WIP)

Front-end

index.html

Simple HTML form

POSTs submitted username to script via Flask

result.html

Another simple HTML template

Rendered from project.py with resulting data

Currently mostly wrestling with Flask

Demo

[Warning: It's not pretty]



Stats for user @hajimesyacho

Percent of JLPT kanji used:

- N5: 87.34%
- N4: 86.75%
- N3: 66.76%
- N2: 34.88%
- N1: 15.58%

Percent of Jōyō kanji used:

- G1: 85.00%
- G2: 75.00%
- G3: 72.50%
- G4: 52.00%
- G5: 43.24%
- G6: 36.46%
- G7: 20.77%

Most common kanji:

[('日', 255), ('今', 214), ('動', 171), ('画', 170), ('寝', 83), ('大', 75), ('気', 73), ('食', 57), ('張', 53), ('人', 48), ('頑', 46), ('出', 42), ('間', 42), ('丈', 40), ('夫', 40), ('岡', 39), ('最', 38), ('様', 38), ('来', 37), ('時', 36)]

Similar writing styles:

WIP...

Stats for user @nhk_news

Percent of JLPT kanji used:

- N5: 100.00%
- N4: 98.80%
- N3: 94.55%
- N2: 86.92%
- N1: 53.49%

Percent of Jōyō kanji used:

- G1: 98.75%
- G2: 95.00%
- G3: 95.50%
- G4: 94.50%
- G5: 92.43%
- G6: 91.16%
- G7: 61.55%

Most common kanji:

[('日', 657), ('大', 592), ('国', 533), ('人', 526), ('会', 496), ('北', 425), ('本', 379), ('朝', 376), ('発', 371), ('鮮', 362), ('年', 341), ('事', 328), ('議', 308), ('中', 295), ('新', 286), ('相', 281), ('見', 262), ('開', 261), ('子', 253), ('場', 245)]

Similar writing styles:

WIP...

Final Notes

Remaining Tasks

— — —

- Stop relying on pre-generated .xml file
- Figure out how to deal with Auth keys/secrets
- Keep trying to get the tokenizer working
- Prettify
- Add WaniKani statistics if there is time

Questions?

Thank you!

(' ▽ ') /