

To clear up:

- lateness to class?
- late on assignments?
- ~~not showing~~ missed class? \rightarrow tell me what's going on.
- missed assignment? { min. completion contract for a unit. (LL - need to assign grades)

To day: question genre: "what makes two types of lang. different?"

altruistic vs. not, funny v. not, ~~in one ti~~ before an event vs. after an event, dem. vs. repub.

- many approaches; today, an example ~~Bayesian~~ statistical approach from a Bayesian perspective (as opposed to a classification perspective)
- \rightarrow goal: inspire you to develop new tests.

Desiderata:

refinement: we want significant differences:

... when diffs ~~are~~ (probably) not due to chance?

(explained. avg $N=60K$)
 \Rightarrow # possible 5-grams = $(60K)^5$
that not enough docs \Rightarrow need a way to talk about things being due to chance

- (a) ~~we'll want to be able to estimate variance of an estimator~~
 - (b) want estimates of the variance of our 'difference' statistics.
 - (c) in lang, generally never enough data. (retated this is a related point)
- and this will certainly be the case for your settings, where you're generally looking @ sthg "tricky" / restricted

b/c of our new domains, not pre-filter features, as far as computationally possible - how do you know what will be distinguishing? ex: un/fair for gender, still "0" for sentiment, 2/4 acant, keep punctuation, #s in Twitter (#s)

Zipf's law (gloss): ~~rarer~~ rare ~~words~~ ~~make up~~ complex surprisingly large fraction of data. \leftarrow prob of x prop. to $\frac{1}{(\text{rank of } x)^2}$. So tail is "heavier" than an exp. decr. \leftarrow <is Zipf's law an example of this?>

idea: use priors to express ~~prob~~ general knowledge about lang, ameliorate the sparse data problem. [even if never seen a dem or repub doc, you know "the" is gonna be relatively frequent]

idea: use priors to express ~~prob~~ general knowledge about lang, ameliorate the sparse data problem. [even if never seen a dem or repub doc, you know "the" is gonna be relatively frequent]

idea: use priors to express ~~prob~~ general knowledge about lang, ameliorate the sparse data problem. [even if never seen a dem or repub doc, you know "the" is gonna be relatively frequent]

\Rightarrow that's when we're gonna be headed \rightarrow distinguished political scientists who also collaborated w/ computer scientists...

main idea of Monroe / ~~Colarusi~~ / ~~Qu~~ et al

"Fightin' Words: Lexical Feature Selection and Evaluation for Identifying the Content of Political Conflict", Political Analysis 2008

contrasts: { logistic regression coefficients (e.g. Mishra; Gilbert) \rightarrow altho these turn out to be log-odds ratios = change in log-odds expected for a one-unit change in corresponding standardized var, other vars held constant. (if no interaction)

statistical { state-based coin-flipping for burstiness detection (Klemm '02)

~~statistical but no explicit prior or estimate of variance.~~

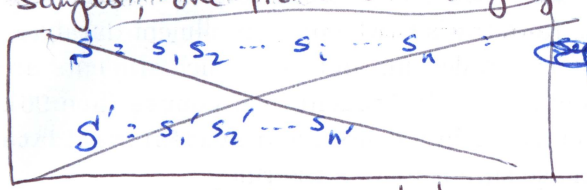
\rightarrow underlying Markov model

ex ref: FAQ how do I interpret odds ratios in logistic regress?

(c) in real research: use held-out set for exploring diffs, if want to be proper and plan to do classification experiments
 (True, it's sad that you lose even more data this way)

Setting: ~~two samples~~ data = lang seg 1: s_1, s_2, \dots, s_n and the other from the other, say Dem. Repub. on topic of abortion

running example



[can generalize to features,] like POS tags, but want some notion of plausible independence has diff. prob of occ.

like to be able to see that, say, a word like 'pro-life' w.r.t. the two ~~under~~ generation processes.

later

Assume ~~words~~ ~~are~~ drawn ~~indep~~, i.i.d. from a multinomial ~~with~~ w. params

$$\theta^D = (\theta_{v_1}^D, \dots, \theta_{|V|}^D)$$

$$\theta^R = (\theta_{v_1}^R, \dots, \theta_{|V|}^R)$$

consider the count histograms: - our observed evidence: $(c_1^D, c_2^D, \dots, c_{|V|}^D)$ vs. $(c_1^R, c_2^R, \dots, c_{|V|}^R)$
 (summary of the) \leftarrow assumed to be a stream - concatenate the docs

following MCQ: 1st gather intuitions, by looking @ simple statistics and establish notation

[3.2.1: diff. in frequencies:] rank v_i by $c_i^D - c_i^R$? \rightarrow give away

- what if one side speaks more?
- be carefully about cherry-picking (p. 376)
- ranking ^{given in paper} makes it more clear ~~the~~: "the", "of", "so"

3.2.2 diff. in proportions: normalize: $\frac{c_i^D}{\sum_j c_j^D} - \frac{c_i^R}{\sum_j c_j^R}$
 rf_i^D (rel. freq.) $\quad rf_i^R$ (rel. freq.)

See Fig 1, p. 377

Great visualization idea by MCQ!



ranking of top 20 D, top 20 R

\Rightarrow most freq words are getting most weight
 most weighted words are all high-freq.

claim: variation in high-freq words not accounted for:

low freq diffs: $\frac{2}{N_0} - \frac{1}{N} = \frac{1}{N}$; high-freq diffs: $\frac{2000}{N} - \frac{1990}{N} = \frac{10}{N} \gg \frac{1}{N}$

(not needed, but build intuition for reinforcement)

b/c we know we're gonna want to get to a log-odds ratio:
 (freshadow: research drama - this won't go well, but is leading to sth good >

let's consider odds ~~ratio~~

$$\frac{r_i^D}{1 - r_i^D}$$

scale is ~~not~~ outside (0,1) - more extreme gives bigger #
 consider 4:1 odds: that the word will appear for Dem

(why are Dem's @??)

log-odds-ratio ~~is~~

$$\log \left(\frac{r_i^D / (1 - r_i^D)}{r_i^R / (1 - r_i^R)} \right)$$

- if equal rel freqs, get 0.
- informative sign.
- problem of for unseen words

see fig 2 - high-wt words are all low-freq, b/c ratios are more extreme for low counts (same 'nooding' % example as before)

let's go back to model-based approach:

→ use notation from previous page that was deferred

- Should we keep trying ad hoc fixes? Or try sth more "principled"?
 let's try to explicitly take variance into acct: if we see a big log-odds ratio, need: what is the variance of the log-odds ratio?

could just smooth by adding in small fake counts to everything, normalize appropriately.
 log-odds ratio is that just by chance?
 why? We then can do a 3-score test: Remember 1.96??

↳ what is its distribution?

model: multinomial generation (via notation that was deferred from previous page)
 - natural, ~~also~~ for independent rolls of dies or draws of words.
 (altho' not true (Noriegar', Ken Church))

what do we know about $\vec{\theta}^D, \vec{\theta}^R$

express prior: Dirichlet is very convenient: conjugate prior for multinomial.
 (not that we know any other priors on multinomials).

parameters $\alpha_1, \alpha_2, \dots, \alpha_{|V|}$; call $\alpha_0 = \sum \alpha_i$; require $\alpha_i > 0$
 (maybe don't need to introduce)

probably better to have just talked about the mean.

prior multinomial where the θ_i 's are roughly normalized

mode: $\theta_i = \frac{\alpha_i - 1}{\sum_j (\alpha_j - 1)}$

- so, like multinomials like this,

if $\alpha_i > 1$ (Albert & Denis, 2012 as one 'random' reference).
 (Note: the "pts" being drawn here represent multinomial "an urn of urns", Matthew said)

described as the "normalization" trick.

variance $\sigma \downarrow$ as α_i 's increase.
 mean: $\theta_i = \frac{\alpha_i}{\alpha_0}$

already guessed use google n-gram corp

showed fig 4 - uninformative prior
 fig 5 - informative prior
 (very definitely made an impression)

posterior $P(\vec{\theta} | \vec{z}) = \text{pmp to } P(\vec{z} | \vec{\theta}) P(\vec{\theta}) \rightarrow$ is Dir w/ $\alpha_i = \alpha_i + c_i$!