# Recommender Systems 2

Internet Analytics (COM-308)

Prof. Matthias Grossglauser
School of Computer and Communication
Sciences

# Overview

- Focus on recommendation of text (prose, tags, …)
- Vector space model
  - Each dimension ~ one term (word)
- TF-IDF metric:
  - Frequency in document makes a word important
  - Frequency in many docs makes a word less important
- Probabilistic model for text classification
  - Naïve Bayes: every word is i.i.d. given class
- Smoothing:
  - Dealing with rare words not seen in training

# Basic idea

- Recommend to user $u$ items similar to the ones he/she liked before
  - Collaborative filtering: similar = liked by people who like the same stuff as $u$
  - Content-based: similar = with similar content features as previously liked items
- What features:
  - Context-dependent
  - Images&music: signal properties (hard); meta-information; tags;...
    - Pandora: music genome project, ~ 400 features
  - Text: easiest & most widespread
    - Prose, tags,...

# Vector space model

- Compact description of a document
  - Ignores order – "bag of words"
- One dimension per term/word
  - Typically very sparse
- Count vector:
  - $f_i$ = # of occurrences of word $i$ in document
- Note: not reversible, ignores order of words
  - The meaning of a sentence would be lost on a human reader!
  - (a a be human lost meaning of on reader sentence the would!)

# Profile from words

- How to create a useful profile of a document?
  - Frequent words are characteristic of "topic"
  - Document A: ("Probability":50, "Markov":20, "Poisson":15,…)
  - Document B: ("Wimbledon":30, "Federer":8, "Nadal":5,…)
- TF: Term Frequency
  - Function of one document $j$ (not the whole corpus)
  - Def: $f_{ij}$ = # of occurrences (frequency) of word $i$ in doc $j$
  - Def: $TF_{ij} = \dfrac{f_{ij}}{\max\limits_{k} f_{kj}}$
  - Importance of word $i$ in document $j$

# TF-IDF: A measure of word importance

- Problem:
  - Most frequent terms would be (in English):
    the, be, to, of, and, a, in, that, have, I, it, for, not, on, with, he, as, you, do, at,...
  - No information, because common to all docs
  - We want words that are frequent **only** in target docs
- IDF: Inverse Document Frequency
  - Function of whole corpus
  - Def: $n_i = \#$ documents $j$ where word $i$ occurs (at least once)
  - Def: $IDF_i = -\log_2 \frac{n_i}{N}$
  - If I know word $i$, number of bits of information I learn about which document is the target within corpus

# TF-IDF vector space model

- Document profile $D$ within a corpus:
  - $TFIDF_{ij} = TF_{ij} \times IDF_i$
  - Take top terms as document profile
  - High score: word frequent in this document, but not in most others
- Vectors are high-dim but sparse
- Refinements: text preprocessing
  - Remove stop words: the, be, to, of, and, a,…
  - Stemming & lemming: transforming
    - "the boy's cars are different colors" -> "the boy car be differ color" [Manning et al.]
  - Vector cutoff to most important terms
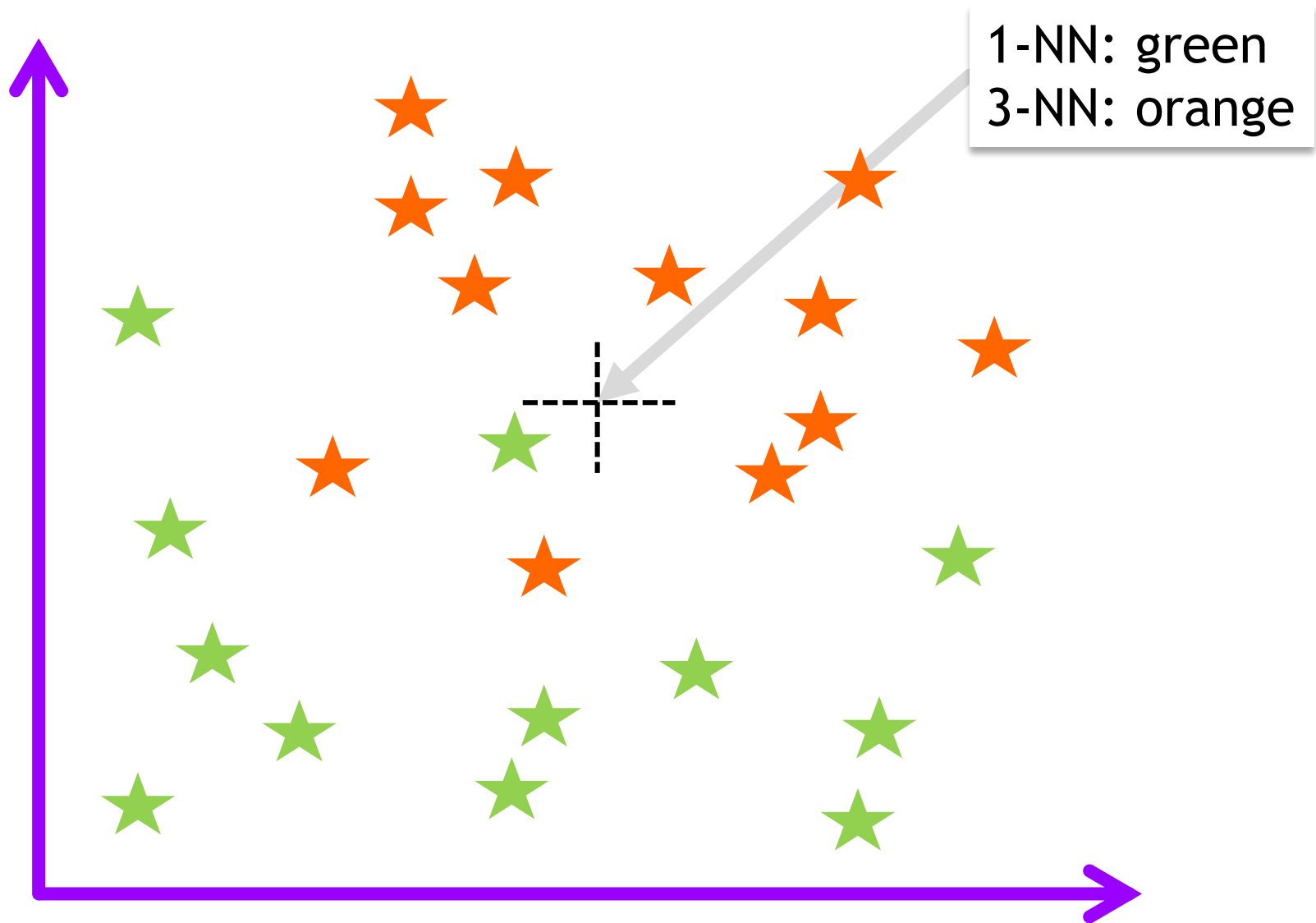  - Allow multi-word terms ("United States")

# Queries and recommendations

- User profile (query) $Q$:
  - Explicit: e.g., formulating a query ("north korea")
  - Implicit: ratings (e.g., "likes")
- Explicit:
  - These models are from information retrieval:
    - Searching by query: return most similar docs to query
    - Query terms → TF-IDF vector $Q$
- Assumption:
  - Likelihood that user $Q$ likes document $D \sim sim(Q, D)$
  - Options for $sim(Q, D)$:
    - $$sim(Q, D) = \cos(Q, D)$$
    - $$sim(Q, D) = \frac{\sum_i (q_i - \bar{q})(d_i - \bar{d})}{\sqrt{\sum (q_i - \bar{q})^2 \sum (d_i - \bar{d})^2}}$$
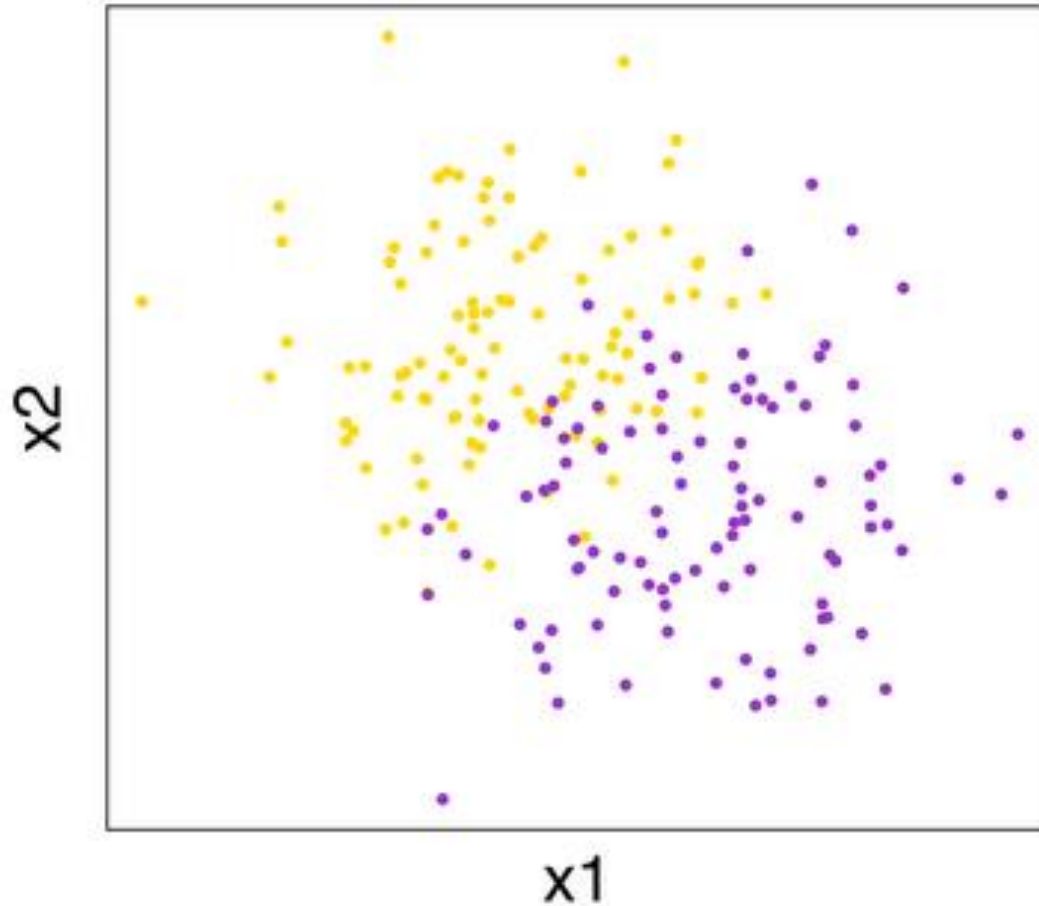
# From queries to ratings

- Implicit: user rates documents rather than queries:
  - Treat highly rated/liked docs as "positive queries", low rated/not liked as "negative queries"
- How to rate a new document $D$?
  - Classification problem: many methods
  - Generic non-parametric method: kNN ($k$ nearest neighbors)
  - Select $k$ rated docs in $Q$ closest to $D$ according to $sim(Q, D)$; majority in this set is predictor

# kNN classifier

1-NN: green
3-NN: orange

# kNN classifier: learning $k$



**Binary kNN Classification Training Set**

[Burton DeWilde: Data Science Rules (datasciencerules.blogspot.com), Oct 2012]
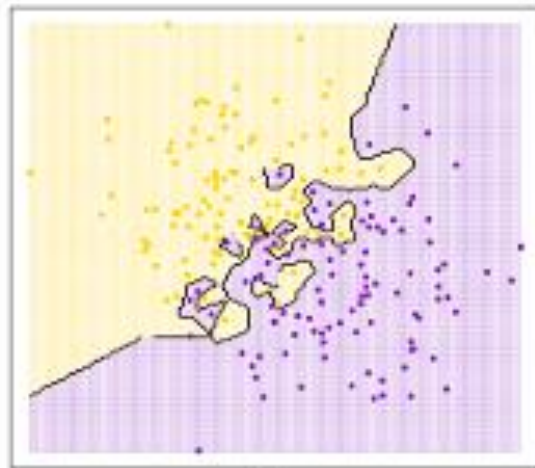
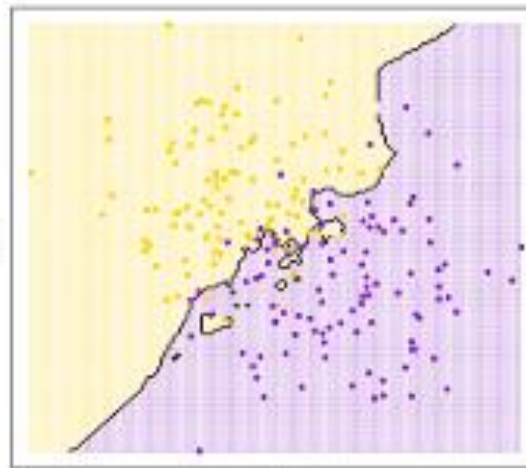# kNN: impact of $k$

overfitted    best model    overgeneralizes



Binary kNN Classification (k=1)    Binary kNN Classification (k=5)    Binary kNN Classification (k=25)

[Burton DeWilde: Data Science Rules (datasciencerules.blogspot.com), Oct 2012]

# Critique of vector-space approach

- Assumptions implicit in approach
    - "small angle between TF-IDF vectors means document close to query": intuitively ok, hard to quantify
    - Quantities do not have "physical meaning", purely heuristic
- We would like a clean model: assumptions, performance measure we can optimize & compare
    - Probabilistic model: rigorous treatment of uncertainty

# Probabilistic models

- Significant uncertainty in predictions
  - Quantization effects: like/dislike -> how much?
  - Context: e.g.: dislike right now (mood), or dislike categorically?
  - Errors, confusions, etc.
- Uncertainty → model explicitly as probability
  - Make assumptions explicit
  - Easier to interpret significance
  - Result comes with measure of uncertainty (confidence interval, etc.)

# ML: supervised vs unsupervised

- ## Supervised learning:
  - Given: (input, output) → find input-output map

$$x_1$$
$$x_2$$
$$x_3$$
→
$$y_1$$
$$y_2$$
$$y_3$$

- ## Unsupervised learning:
  - Given: (input) → find structure

$$x_1$$
$$x_2$$
$$x_3$$
structure

# Supervised ML: classification vs regression

- ## Classification:
  - *Y* is a class label
  - Categorical (not necessarily numeric)
  - Example: {spam, not spam}; {blue, green}; political party affiliation inferred from questions
- ## Regression:
  - *Y* is (typically) in $\mathbb{R}$
  - Example: temperature tomorrow; ranking for a movie (1...5 stars)

# Bayesian inference

- Statistical inference: frequentist (non-Bayesian)
  - Observation $Z$
  - Model: $p_\theta(z)$: distribution of $Z$, depending on hidden parameter $\theta$
  - Goal: infer $\theta$ from observation(s) of $Z$
  - Maximum Likelihood estimator: $\hat{\theta} = \max_\theta p_\theta(Z)$
    - Estimated parameter best explains observed data

# Bayesian inference

- Statistical inference: Bayesian
  - We know something about $\theta$: prior knowledge about the problem
  - $\theta$ is a random variable with a known distribution: prior
  - Model: $p(Z|\theta)$: distribution of $Z$, conditional on hidden random variable $\theta$
  - Bayes' rule:
  $$P(\theta|Z) = \frac{P(\theta, Z)}{P(Z)} = \frac{P(Z|\theta)P(\theta)}{\sum_{\theta'} P(Z|\theta')P(\theta')}$$
  - Maximum A Posteriori (MAP) estimator:
  $$\hat{\theta} = \max_{\theta} P(\theta|Z)$$
  - But the full posterior distribution $P(\theta|Z)$ carries additional information!
    - How certain/uncertain are we about $\theta$ given data $Z$

# Example: Max-Likelihood vs Bayesian

- Medical test
  - You take a medical test whose accuracy is 90% - that is, prob. test gives right result = 0.9
  - Frequentist:
    - $P(pos|sick) = 0.9$; $P(pos|healthy) = 0.1$
    - ML: $Z = pos \rightarrow \hat{\theta} = sick$
    - Test comes back positive → you conclude you are sick

# Example: ML vs Bayesian

- Medical test:
  - Bayesian:
    - Medical test; prior = one in a million: $P(sick) = 10^{-6}$
    - If test comes back positive:
    - $P(sick|pos) = \dfrac{P(pos|sick)P(sick)}{P(pos|sick)P(sick)+P(pos|healthy)P(healthy)}$
    - $P(sick|pos) \cong 0.9 \times 10^{-5}$
    - You conclude you are very likely healthy!
  - Watch out: doctors apparently don't know this!

# Naïve Bayes classifier

- Need a probabilistic model for a document
- Simplest model:
  - Naïve = independent terms (features)
  - Each word is generated according to i.i.d. distribution

$$P(Z_1, \ldots, Z_n | \theta) = \prod_i P(Z_i | \theta)$$

- Hidden variable:
  - Relevant (good, $G$) or not relevant (bad, $B$)
- Observable variable:
  - Message = set of words $(z_1, z_2, \ldots, z_n)$
- Classify message into $(G, B)$
- Model $p(Z|\{G, B\}), p(\{G, B\})$:
  - Learn from data

# Example: naïve Bayes classifier

- Training set:

| Get nice watch |
|---|
| New York rocks! |
| Watch for rocks |

| Cheap replica watch |
|---|
| New cheap loan |
| Get lottery million |
| Million dollar watch |

- Prior: $P(\theta = G) = \frac{3}{7}$; $P(\theta = B) = \frac{4}{7}$

- Conditional word distributions $P(Z|\theta)$:

| Z | get | nice | watch | new | york | rocks | for | cheap | replica | loan | lottery | million | dollar | perfect |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $9 \times P(Z\|G)$ | 1 | 1 | 2 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $12 \times P(Z\|B)$ | 1 | 0 | 2 | 1 | 0 | 0 | 0 | 2 | 1 | 1 | 1 | 2 | 1 | 0 |

# Example: naïve Bayes classifier

- Classifying sentences $M = (Z_1, Z_2, Z_3, \ldots)$:
  - «get new watch»:

$$P(G|M) =$$

$$= \frac{P(Z_1|G)P(Z_2|G)P(Z_3|G)P(G)}{P(Z_1|G)P(Z_2|G)P(Z_3|G)P(G) + P(Z_1|B)P(Z_2|B)P(Z_3|B)P(B)} =$$

$$= \frac{9^{-3} \cdot 1 \cdot 1 \cdot 2 \cdot 3/7}{9^{-3} \cdot 1 \cdot 1 \cdot 2 \cdot \frac{3}{7} + 12^{-3} \cdot 1 \cdot 1 \cdot 2 \cdot 4/7} = 0.64$$

| Z | get | nice | watch | new | york | rocks | for | cheap | replica | loan | lottery | million | dollar | perfect |
|---|-----|------|-------|-----|------|-------|-----|-------|---------|------|---------|---------|--------|---------|
| 9 $\times P(Z|G)$ | 1 | 1 | 2 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 $\times P(Z|B)$ | 1 | 0 | 2 | 1 | 0 | 0 | 0 | 2 | 1 | 1 | 1 | 2 | 1 | 0 |

# Example: naïve Bayes classifier

- Classifying sentences $M = (Z_1, Z_2, Z_3, \ldots)$:
  - «cheap replica rocks»:

$$P(G|M) =$$

$$= \frac{\boxed{=0} \quad P(Z_3|G)P(G)}{\boxed{=0} \quad P(Z_3|G)P(G) + P(Z_1|B)P(Z_2|B) \boxed{=0} \quad P(B)}$$

  - Undefined!

| Z | get | nice | watch | new | york | rocks | for | cheap | replica | loan | lottery | million | dollar | perfect |
|---|-----|------|-------|-----|------|-------|-----|-------|---------|------|---------|---------|--------|---------|
| $9 \times P(Z|G)$ | 1 | 1 | 2 | 1 | 1 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| $12 \times P(Z|B)$ | 1 | 0 | 2 | 1 | 0 | 0 | 0 | 2 | 1 | 1 | 1 | 2 | 1 | 0 |

# Problem with unseen training terms

- ## Sparsity problem:
  - If alphabet of words is large w.r.t. training set, there are some words $z$ we never see (e.g., $z =$ "mesonoxian")
    - Estimate: $P(\text{mesonoxian}|\{G,B\}) = 0$
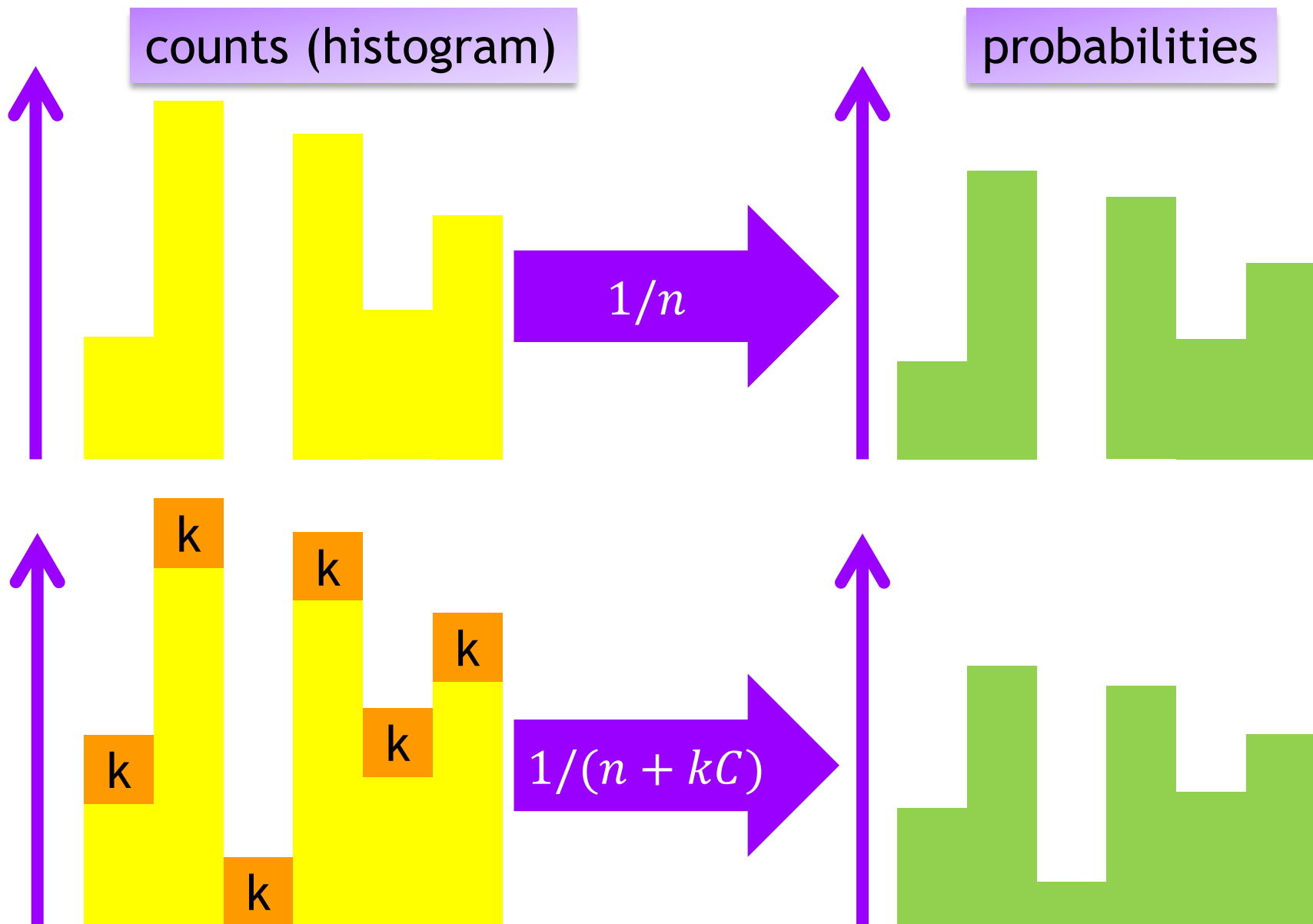  - If target message contains "mesonoxian":
    $$P(\{G,B\}) = \frac{P(z|\theta)P(\theta)}{\sum_{\theta'} P(z|\theta')P(\theta')} = \frac{0}{0}$$

- ## Problem:
  - We estimate a distribution from a very small set of samples – a form of overfitting
  - How to correctly estimate very rare words?

- ## Approach 1:
  - Ignore unseen words → simple, but crude; throws away information

# Laplace smoothing

- Idea: assume every word occurs at least once
  - Aka "additive smoothing", "add-one smoothing"
- Bias towards uniform distribution
  - A form of regularization
- Estimate of a distribution over domain $D = \{d_1, \ldots, d_C\}$ from data set $\{x_1, x_2, \ldots, x_n\}$
  - Unsmoothed: $p(X = x) = \dfrac{|\{x : x = d_i\}|}{n}$ ($n$=# samples)
  - Smoothed: assume $k$ "fake" observations for each class
  $$p(X = d_i) = \frac{|\{x : x = d_i\}| + k}{n + kC}$$
  - Empty dataset ($n = 0$) → $P(Z|\theta)$ uniform
  - Large dataset ($n \gg 1$) → smoothed $P(Z|\theta) \cong$ unsmoothed $P(Z|\theta)$

# Laplace smoothing

counts (histogram)

probabilities

$1/n$

k

k

k

k

k

k

$1/(n + kC)$

# Example: Laplace-smoothed classifier

- Sentence $M =$«cheap replica rocks»:

$$P(G|M) =$$

$$= \frac{P(Z_1|G)P(Z_2|G)P(Z_3|G)P(G)}{P(Z_1|G)P(Z_2|G)P(Z_3|G)P(G) + P(Z_1|B)P(Z_2|B)P(Z_3|B)P(B)} =$$

$$= \frac{23^{-3} \cdot 1 \cdot 1 \cdot 3 \cdot 4/9}{23^{-3} \cdot 1 \cdot 1 \cdot 3 \cdot 4/9 + 26^{-3} \cdot 3 \cdot 2 \cdot 3 \cdot 5/9} = 0.37$$

- Advantages:
  - We can compute an estimate for any message
  - For small training sets → avoids overfitting

| $Z$ | get | nice | watch | new | york | rocks | for | cheap | replica | loan | lottery | million | dollar | perfect |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $23 \times P(Z|G)$ | 2 | 2 | 3 | 2 | 2 | 3 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| $26 \times P(Z|B)$ | 2 | 1 | 3 | 2 | 1 | 1 | 1 | 3 | 2 | 2 | 2 | 3 | 2 | 1 |

# RecSys: content vs collaborative

| Pros |
|---|
| Independent of other users → no cold start problem for new items (item comes with features) |
| Independent of other users → can recommend for unique tastes, no "trend to average" |
| Can provide reasons for recommendation (e.g., matching keywords) |

| Cons |
|---|
| Multimedia etc.: hard to identify features |
| Independent of other users → no discovery, no surprises |
| Cold start problem for new user |

- In practice: combination
  - Lack of ratings, few users → rely more on content
  - Lots of users, few tags → collaborative

# Summary

- Content: text, tags, user comments, subtitles,…
- Collaborative filtering vs content-based:
  - Blind to content vs blind to other users
- Classical approaches from information retrieval:
  - Vector space models, similarity metrics
- Modern probabilistic approaches from ML:
  - Naïve Bayes, language models ($n$-grams), word embeddings
- Other application for naïve Bayes: spam filtering
  - $P(B) \cong 0.8 \dots 0.9$

# References

- [A. Rajaraman, J. D. Ullman: Mining of Massive Datasets, Cambridge, 2012 (chapter 9)]

- [S. Russell, P. Norvig: Artificial Intelligence – A Modern Approach (3$^{rd}$ ed), Pearson, 2010 (chapter22)]

- [W. B. Croft, D. Metzler, T. Strohman: Search Engines – Information Retrieval in Practice, Addison Wesley, 2010 (chapters 7&10)]

- [Ch. D. Manning, P. Raghavan, H. Schütze: Introduction to Information Retrieval, Cambridge, 2008]