

NLP: Classification

AMITDHOMNE

1 Classification Tasks

- Language identification: determine the language that a text is written in
- Spam filtering: label emails, tweets, blog comments as spam or not spam
- Routing: label emails to appropriate people in an organization (complaints, tech support, order status, etc)
- Sentiment analysis: label some text as being positive or negative (polarity classification)

Example: Sentiment

- Determine the sentiment (positive vs. negative) of, for example, a tweet
- “Probably the worst movie of the century”
- One idea: Compare of “positive” words (good, great, best) to “negative” words (bad, terrible, worst).
 - Humans give insufficient lists. Learning a sufficient list is hard.
 - Words mean different things in different contexts:
 - “This thing is *a great deal*. Definitely worth the money.”
 - “*A great deal* of media attention surrounded the event.”
 - “It’s *a great deal*... if you’re looking to looking to waste your money.”
 - Some words can flip polarity:
 - “It’s **not** a **good** investment.”
 - “I thought it would be **terrible**, **but** I was so wrong.”
 - Multi-word expressions:
 - “The movie was shit.”
 - “The movie was the shit.”
 - Subtlety:
 - “If this movie’s your thing, don’t bother talking to me.”
 - “Great plot, great acting, great cast, but it just doesn’t hold up.”

- Can depend on the target:
“Unpredictable plot”
“Unpredictable steering”
- Additional Features
 - Ngrams: “must buy”, “couldn’t care less”
 - Casing: uppercase words are often subjective
 - Punctuation: lots of ! or ? can indicate subjectivity
 - Emoticons: :) vs. :(

2 Rule-Based System

- Write prediction rules: “If contains X and Y but not Z, then ‘positive’”
- What happens when multiple rules apply but conflict?
 - Order the rules according to their accuracy?
 - Assign weights to the rules?
- Problems
 - Time-consuming and expensive to write rules.
 - High precision, low recall
 - Rules have to be manually tailored to each dataset (unpredictable vs. unpredictable)
 - Expensive to update (new expressions, new slang, new abrvs)

3 Learning

- If we have examples, we can learn a function mapping texts to categories
- Often probabilistic
- Instead of rules, we use features
- Features are automatically weighted based on statistics in the training data.
- Features are dimensions in space.
- Learn a boundary between classes.
- Boundary used to classify new texts.

Some Data

```

start=B, end=ia, location
start=B, end=er, person
start=M, end=ia, person
start=L, end=ia, location
start=N, end=er, location
start=B, end=ia, location
start=E, end=nd, location
start=N, end=ia, location
start=A, end=er, person
start=L, end=ke, person

```

Our Goal

- Learn a function that maps features to the most likely label
- $\text{best_label} = \text{argmax}_{\text{label}} p(\text{label} \mid \text{features})$

Direct Posterior Parameter Estimation from Data

- $p(\text{label} \mid \text{features}) = \frac{C(\text{instances with label and feature})}{C(\text{instances with features})}$
 - $p(\text{label} = \text{location} \mid \text{start} = B, \text{end} = \text{er}) = \frac{0}{1} = 0.0$
 - $p(\text{label} = \text{person} \mid \text{start} = B, \text{end} = \text{er}) = \frac{1}{1} = 1.0$
 - $p(\text{label} = \text{location} \mid \text{start} = B, \text{end} = \text{ia}) = \frac{2}{2} = 1.0$
 - $p(\text{label} = \text{person} \mid \text{start} = B, \text{end} = \text{ia}) = \frac{0}{2} = 0.0$
- This doesn't work very well
- Sparsity: any particular feature combination is rare, hard to generalize
- Even worse when every word in a text is a feature
 - Every text would be unique, and there would be no generalization at all
 - Couldn't label new instances

Bayes Rule

- $p(\text{label} \mid \text{features}) = \frac{p(\text{features} \mid \text{label}) \cdot p(\text{label})}{p(\text{features})}$
- $\text{best_label} = \text{argmax}_{\text{label}} \frac{p(\text{features} \mid \text{label}) \cdot p(\text{label})}{p(\text{features})}$
- If labels = {A,B}: $\frac{p(\text{features} \mid A) \cdot p(A)}{p(\text{features})}$ vs. $\frac{p(\text{features} \mid B) \cdot p(B)}{p(\text{features})}$
- Denominator is always the same, so: $p(\text{features} \mid A) \cdot p(A)$ vs. $p(\text{features} \mid B) \cdot p(B)$
- Thus, to compute the **posterior**, we need two things
 - the likelihood of the evidence: $p(\text{features} \mid \text{label})$

- the prior: $p(label)$

Direct Evidence Likelihood Estimation from Data

- $p(features \mid label) = \frac{C(instances \text{ with } features \text{ and } label)}{C(instances \text{ with } label)}$
 - $p(start = B, end = er \mid label = location) = \frac{0}{6} = 0.0$
 - $p(start = B, end = er \mid label = person) = \frac{1}{4} = 0.25$
 - $p(start = B, end = ia \mid label = location) = \frac{2}{6} = 0.33$
 - $p(start = B, end = ia \mid p(label = person) = \frac{0}{4} = 0.0$
- Still problematic: still sparse, still lots of zeros, still hard to generalize

Naïve Bayes

- We want to disentangle the features for better generalization
- Compute each feature's probability independently
- Will be able to compute the probability of an instance from the features even if we haven't seen that particular combination of features before.
- Requires us to assume that features are independent
 - Not actually true! Language doesn't work like that.
 - But it's a simplifying assumption
 - “Naïve” assumption
- $p(features \mid label) = p(F_1, F_2, F_3, \dots \mid label) = p(F_1 \mid label) \cdot p(F_2 \mid label) \cdot p(F_3 \mid label) \cdot \dots$

Parameter Estimation from Data

- The prior
 - $p(label = location) = \frac{C(label=location)}{\sum_l C(label=l)} = \frac{6}{10} = 0.6$
 - $p(label = person) = \frac{C(label=person)}{\sum_l C(label=l)} = \frac{4}{10} = 0.4$
- Likelihood of the evidence
 - $p(start = B \mid label = location) = \frac{C(start=B, label=location)}{C(label=location)} = \frac{2}{6} = 0.33$
 - $p(start = B \mid label = person) = \frac{C(start=B, label=person)}{C(label=person)} = \frac{1}{4} = 0.25$
 - $p(end = ia \mid label = location) = \frac{C(end=ia, label=location)}{C(label=location)} = \frac{4}{6} = 0.67$
 - $p(end = ia \mid label = person) = \frac{C(end=ia, label=person)}{C(label=person)} = \frac{1}{4} = 0.25$

$$\begin{aligned}
- p(\text{end} = \text{nd} \mid \text{label} = \text{location}) &= \frac{C(\text{end}=\text{nd}, \text{label}=\text{location})}{C(\text{label}=\text{location})} = \frac{1}{6} = 0.17 \\
- p(\text{end} = \text{nd} \mid \text{label} = \text{person}) &= \frac{C(\text{end}=\text{nd}, \text{label}=\text{person})}{C(\text{label}=\text{person})} = \frac{0}{4} = 0.0
\end{aligned}$$

Naïve Probabilities

- Before

$$\begin{aligned}
- p(\text{start} = B, \text{end} = \text{ia} \mid \text{label} = \text{location}) &= \frac{2}{6} = 0.33 \\
- p(\text{start} = B, \text{end} = \text{ia} \mid p(\text{label} = \text{person})) &= \frac{0}{4} = 0.0
\end{aligned}$$

- Now

$$\begin{aligned}
- p(\text{start} = B \mid \text{label} = \text{location}) \cdot p(\text{end} = \text{ia} \mid \text{label} = \text{location}) &= 0.33 \cdot 0.67 = 0.22 \\
- p(\text{start} = B \mid \text{label} = \text{location}) \cdot p(\text{end} = \text{ia} \mid \text{label} = \text{person}) &= 0.25 \cdot 0.25 = 0.06
\end{aligned}$$

Classifying

- We get a **new** instance. Need to determine its label.
- Works even if we haven't seen the particular combination of features
- **start=L, end=er**
- $p(\text{features} \mid A) \cdot p(A)$ vs. $p(\text{features} \mid B) \cdot p(B)$
- $p(\text{start} = L \mid \text{location}) \cdot p(\text{end} = \text{er} \mid \text{location}) \cdot p(\text{location}) = \frac{1}{6} \cdot \frac{1}{6} \cdot \frac{6}{10} = 0.02$
- $p(\text{start} = L \mid \text{person}) \cdot p(\text{end} = \text{er} \mid \text{person}) \cdot p(\text{person}) = \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{4}{10} = 0.06$
- So, it's more likely a *person*

Why Naïve Bayes?

- Modularity: separate out individual features and prior
- Helps us deal with sparsity
 - Particular feature combinations are rare
 - Individual features are less sparse
 - Still have some features that appear only with one label (meaning zero probabilities), but this is less common
- Priors
 - Controls how much the base label distribution affects the probabilities
 - Can be automatically calculated from data (as we've seen)
 - Can be set from outside knowledge, if available
 - * Imagine we are explicitly told that 3/4 of named entities are people
 - * $p(\text{label} = \text{person}) = 0.75$

- * $p(\text{start} = L \mid \text{location}) \cdot p(\text{end} = er \mid \text{location}) \cdot p(\text{location}) = \frac{1}{6} \cdot \frac{1}{6} \cdot 0.25 = 0.007$
- * $p(\text{start} = L \mid \text{person}) \cdot p(\text{end} = er \mid \text{person}) \cdot p(\text{person}) = \frac{1}{4} \cdot \frac{2}{4} \cdot 0.75 = 0.09$
- * So the likelihood of *person* is even higher
- Useful for injecting linguistic knowledge into a learned model

// Day 2

4 Smoothing

Classify *Austria*:

- The prior

$$\begin{aligned}
 - p(\text{label} = \text{location}) &= \frac{C(\text{label}=\text{location})}{\sum_l C(\text{label}=l)} = \frac{6}{10} = 0.6 \\
 - p(\text{label} = \text{person}) &= \frac{C(\text{label}=\text{person})}{\sum_l C(\text{label}=l)} = \frac{4}{10} = 0.4
 \end{aligned}$$

- Likelihood of the evidence

$$\begin{aligned}
 - p(\text{end} = ia \mid \text{label} = \text{location}) &= \frac{C(\text{end}=ia, \text{label}=\text{location})}{C(\text{label}=\text{location})} = \frac{4}{6} = 0.67 \\
 - p(\text{end} = ia \mid \text{label} = \text{person}) &= \frac{C(\text{end}=ia, \text{label}=\text{person})}{C(\text{label}=\text{person})} = \frac{1}{4} = 0.25 \\
 - p(\text{start} = A \mid \text{label} = \text{location}) &= \frac{C(\text{start}=A, \text{label}=\text{location})}{C(\text{label}=\text{location})} = \frac{0}{6} = 0.00 \\
 - p(\text{start} = A \mid \text{label} = \text{person}) &= \frac{C(\text{start}=A, \text{label}=\text{person})}{C(\text{label}=\text{person})} = \frac{1}{4} = 0.25
 \end{aligned}$$

- The posterior

$$\begin{aligned}
 - p(\text{label} = \text{location} \mid \text{start} = A, \text{end} = ia) \\
 &= p(\text{label} = \text{location}) \cdot p(\text{start} = A \mid \text{label} = \text{location}) \cdot p(\text{end} = ia \mid \text{label} = \text{location}) \\
 &= 0.6 \cdot 0.0 \cdot 0.67 = \mathbf{0.0} \\
 - p(\text{label} = \text{person} \mid \text{start} = A, \text{end} = ia) \\
 &= p(\text{label} = \text{person}) \cdot p(\text{start} = A \mid \text{label} = \text{person}) \cdot p(\text{end} = ia \mid \text{label} = \text{person}) \\
 &= 0.4 \cdot 0.25 \cdot 0.25 = \mathbf{0.025}
 \end{aligned}$$

- Since $0.025 > 0.0$, we pick the label *person*
- This seems wrong
 - *end* = *ia* leans strongly toward *location* (0.67 vs. 0.25)
 - *location* is more likely *a priori* (0.6 vs 0.4)
 - While *start* = *A* leans toward *person*, it's not so strong (0.25 vs. 0.00)
 - Moreover, we have very little evidence regarding *start* = *A*, so we probably shouldn't give much weight to it during classification
- Why did we get the wrong answer?
 - Since we had no evidence for *label* = *location* AND *start* = *A*, we got a zero count, which gave us a zero probability.

- So, basically, the $start = A$ feature ended up being the only thing that mattered.
- Even though $start = A$ should have been the least important thing in the calculation.

Add- λ smoothing

- We need to avoid having any zero counts.
- More generally, we need to deal with cases in which we have so few counts that it's hard to generalize
 - $C(A) = 1000, C(B) = 2$ wouldn't be considered a lack of evidence
 - $C(A) = 2, C(B) = 1$ probably doesn't have enough evidence to judge well
- One way to make all the numbers non-zero is to simply add 1 to all of them.
- So if the count was 0, it becomes 1, 1 becomes 2, 100 becomes 101.
- $start$ values = {B,M,L,N,E,A}, end values = {ia,er,nd,ke}

(Feature, Label)	$C(F,L)$	+1	$p(F L)$	$p_{+1}(F L)$	(Feature, Label)	$C(F,L)$	+1	$p(F L)$	$p_{+1}(F L)$
(B, loc)	2	3	$2/6=0.33$	$3/12=0.25$	(B, per)	1	2	$1/4=0.25$	$2/10=0.20$
(M, loc)	0	1	$0/6=0.00$	$1/12=0.08$	(M, per)	1	2	$1/4=0.25$	$2/10=0.20$
(L, loc)	1	2	$1/6=0.17$	$2/12=0.17$	(L, per)	1	2	$1/4=0.25$	$2/10=0.20$
(N, loc)	2	3	$2/6=0.33$	$3/12=0.25$	(N, per)	0	1	$0/4=0.00$	$1/10=0.10$
(E, loc)	1	2	$1/6=0.17$	$2/12=0.17$	(E, per)	0	1	$0/4=0.00$	$1/10=0.10$
(A, loc)	0	1	$0/6=0.00$	$1/12=0.08$	(A, per)	1	2	$1/4=0.25$	$2/10=0.20$
(·, loc)	6	12			(·, per)	4	10		
(ia, loc)	4	5	$4/6=0.67$	$5/10=0.50$	(ia, per)	1	2	$1/4=0.25$	$2/8=0.25$
(er, loc)	1	2	$1/6=0.17$	$2/10=0.20$	(er, per)	2	3	$2/4=0.50$	$3/8=0.38$
(nd, loc)	1	2	$1/6=0.17$	$2/10=0.20$	(nd, per)	0	1	$0/4=0.00$	$1/8=0.13$
(ke, loc)	0	1	$0/6=0.00$	$1/10=0.10$	(ke, per)	1	2	$1/4=0.25$	$2/8=0.25$
(·, loc)	6	10			(·, per)	4	8		

- The count of each (Feature, Label) pair goes up by 1, but the total count for the label across all features goes up by $(1 \cdot (\# \text{ of Feature Values}))$
 - Since we add 1 additional count for each (F,L) — even the ones we've never seen before — and there is one (F,L) pair for each possible F, there are $1 \cdot |\text{FeatureValues}|$ additional counts for the label.
 - $\sum_{f \in \text{Features}} (C(f, L) + 1) = \left(\sum_{f \in \text{FeatureValues}} C(f, L) \right) + |\text{FeatureValues}|$
 - So now, $p(\text{feature} \mid \text{label}) = \frac{C(\text{feature}, \text{label}) + 1}{C(\text{label}) + (1 \cdot |\text{FeatureValues}|)}$
 - If we didn't increase the denominator:

$$p(\text{start} = B \mid \text{label} = \text{location}) = \frac{C(\text{start}=B, \text{label}=\text{location}) + 1}{C(\text{label}=\text{location})} = \frac{2+1}{6} = \frac{3}{6} = 0.50$$

$$p(\text{start} = M \mid \text{label} = \text{location}) = \frac{C(\text{start}=M, \text{label}=\text{location}) + 1}{C(\text{label}=\text{location})} = \frac{0+1}{6} = \frac{1}{6} = 0.17$$

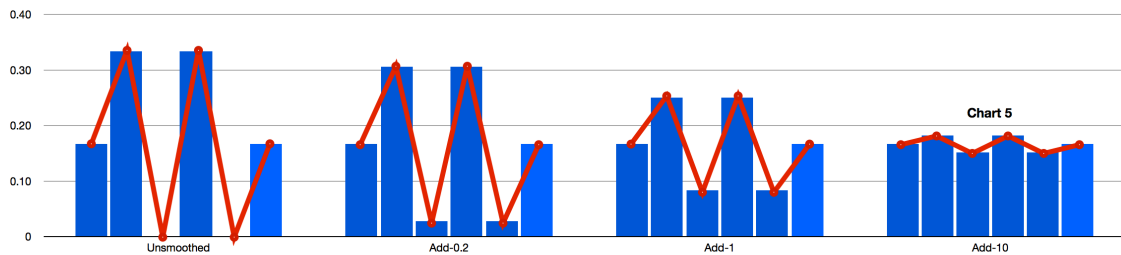
$$\begin{aligned}
p(\text{start} = L \mid \text{label} = \text{location}) &= \frac{C(\text{start}=L, \text{label}=\text{location})+1}{C(\text{label}=\text{location})} = \frac{1+1}{6} = \frac{2}{6} = 0.33 \\
p(\text{start} = N \mid \text{label} = \text{location}) &= \frac{C(\text{start}=N, \text{label}=\text{location})+1}{C(\text{label}=\text{location})} = \frac{2+1}{6} = \frac{3}{6} = 0.50 \\
p(\text{start} = E \mid \text{label} = \text{location}) &= \frac{C(\text{start}=E, \text{label}=\text{location})+1}{C(\text{label}=\text{location})} = \frac{1+1}{6} = \frac{2}{6} = 0.33 \\
p(\text{start} = A \mid \text{label} = \text{location}) &= \frac{C(\text{start}=A, \text{label}=\text{location})+1}{C(\text{label}=\text{location})} = \frac{0+1}{6} = \frac{1}{6} = 0.17
\end{aligned}$$

* But $0.50 \cdot 0.17 \cdot 0.33 \cdot 0.50 \cdot 0.33 \cdot 0.17 = 2.0$

* Since all the values in a probability distribution must sum to 1.0, we no longer have a probability distribution

* Adding the correct amount to the denominator fixes this problem.

- Smoothing flattens out the distribution



- (Partial) counts are “stolen” from high-count items and “given” to low-count items
- Extreme values are smoothed the most
- The impact of smoothing is depends on the relative values of λ and the counts.
- If there are few overall counts, then a relatively large λ will over-smooth, making all values approximately the same likelihood.
- If there are many overall counts, then a relatively small λ will not make much impact.
- $\lambda = 0$ means no smoothing
- $\lambda = \infty$ means uniform distribution (all have equal probability)

- Smoothing changes the relative probabilities of things most when there are low overall counts

- This is good because low counts means low evidence and, thus, low confidence
- We want the impact of evidence to be strong only when we have high confidence
- For naïve Bayes, we always compare $p(F \mid L_1)$ to $p(F \mid L_2)$
- If counts $C(F, L_1)$ and $C(F, L_2)$ are small, we want smoothing to push the relative probabilities together, thus giving them a smaller impact in the NB calculation. (The most extreme example is equal probabilities, which contribute no information to NB vs zero counts, which completely dominate the NB calculation).
- If counts $C(F, L_1)$ and $C(F, L_2)$ are large, we do not want smoothing to have a large impact on the relative probabilities
- For things with low counts (0 vs. 1), smoothing pushes the relative probabilities close together
 - If $C(A) = 1, C(B) = 0$
 - * $p(A) = \frac{1}{1} = \mathbf{1.0}$ and $p(B) = \mathbf{0.0}$

- * $p_{+1}(A) = \frac{2}{3} = \mathbf{0.67}$ and $p_{+1}(B) = \mathbf{0.33}$
- If $C(A) = 2, C(B) = 1$
 - * $p(A) = \frac{2}{3} = \mathbf{0.67}$ and $p(B) = \mathbf{0.33}$
 - * $p_{+1}(A) = \frac{3}{5} = \mathbf{0.60}$ and $p_{+1}(B) = \mathbf{0.40}$
- For things with non-trivial counts, the smoothed relative probabilities don't change much. If $C(A) = 30, C(B) = 0$
 - * $p(A) = \frac{30}{30} = \mathbf{1.0}$ and $p(B) = \mathbf{0.0}$
 - * $p_{+1}(A) = \frac{31}{32} = \mathbf{0.97}$ and $p_{+1}(B) = \mathbf{0.03}$
- If $C(A) = 20, C(B) = 10$
 - * $p(A) = \frac{20}{30} = \mathbf{0.67}$ and $p(B) = \mathbf{0.33}$
 - * $p_{+1}(A) = \frac{21}{32} = \mathbf{0.66}$ and $p_{+1}(B) = \mathbf{0.44}$

Classify *Austria* with Add-1 smoothing

- The prior
 - We do *not* need (or want) to smooth the prior distribution $p(\text{Label})$
 - Smoothing is useful when we have very low or zero counts
 - Counts of labels will never be so low that we can't estimate them
- Likelihood of the evidence
 - $p(\text{end} = ia \mid \text{label} = \text{location}) = \frac{C(\text{end}=ia, \text{label}=\text{location})+1}{C(\text{label}=\text{location})+|\text{end values}|} = \frac{4+1}{6+4} = \frac{5}{10} = 0.50$
 - $p(\text{end} = ia \mid \text{label} = \text{person}) = \frac{C(\text{end}=ia, \text{label}=\text{person})}{C(\text{label}=\text{person})+|\text{end values}|} = \frac{1+1}{4+4} = \frac{2}{8} = 0.25$
 - $p(\text{start} = A \mid \text{label} = \text{location}) = \frac{C(\text{start}=A, \text{label}=\text{location})}{C(\text{label}=\text{location})+|\text{start values}|} = \frac{0+1}{6+6} = \frac{1}{12} = 0.08$
 - $p(\text{start} = A \mid \text{label} = \text{person}) = \frac{C(\text{start}=A, \text{label}=\text{person})}{C(\text{label}=\text{person})+|\text{start values}|} = \frac{1+1}{4+6} = \frac{2}{10} = 0.20$
 - No longer have zeros (this is essential)
 - $\text{start} = A$ relative probabilities are pushed closer together than $\text{end} = ia$ since there is less evidence for either label
- The posterior
 - $p(\text{label} = \text{location} \mid \text{start} = A, \text{end} = ia)$
 $= p(\text{label} = \text{location}) \cdot p(\text{start} = A \mid \text{label} = \text{location}) \cdot p(\text{end} = ia \mid \text{label} = \text{location})$
 $= 0.6 \cdot 0.08 \cdot 0.50 = \mathbf{0.024}$
 - $p(\text{label} = \text{person} \mid \text{start} = A, \text{end} = ia)$
 $= p(\text{label} = \text{person}) \cdot p(\text{start} = A \mid \text{label} = \text{person}) \cdot p(\text{end} = ia \mid \text{label} = \text{person})$
 $= 0.4 \cdot 0.20 \cdot 0.25 = \mathbf{0.020}$
- Since $0.024 > 0.020$, we pick the label *location*
- So smoothing fixed this classification example!

5 Smoothing over unknowable things

- In the tennis example, we have four features, and all of their possible values are known.
- What happens if the set of possible values is unknown?
- What if we want to classify the word “Zambia” as a person or location?

$$- p_{+1}(end = ia \mid location) = \frac{C(end=ia|location)+1}{C(location)+|end\ values|} = \frac{4+1}{6+4} = \frac{5}{10} = 0.50$$

$$- p_{+1}(end = ia \mid person) = \frac{C(end=ia|person)+1}{C(person)+|end\ values|} = \frac{1+1}{4+4} = \frac{2}{8} = 0.25$$

$$- p_{+1}(start = Z \mid location) = \frac{C(start=Z|location)}{C(location)+|start\ values|} = \frac{0+1}{6+6} = \frac{1}{12} = 0.08$$

$$- p_{+1}(start = Z \mid person) = \frac{C(start=Z|person)}{C(person)+|start\ values|} = \frac{0+1}{4+6} = \frac{1}{10} = 0.10$$

- So where does $start = Z$ fit in to the PD?

$$- p_{+1}(start = B \mid location) = \frac{C(start=B|location)}{C(location)+|start\ values|} = \frac{2+1}{6+6} = \frac{3}{12} = 0.25$$

$$- p_{+1}(start = M \mid location) = \frac{C(start=M|location)}{C(location)+|start\ values|} = \frac{0+1}{6+6} = \frac{1}{12} = 0.08$$

$$- p_{+1}(start = L \mid location) = \frac{C(start=L|location)}{C(location)+|start\ values|} = \frac{1+1}{6+6} = \frac{2}{12} = 0.17$$

$$- p_{+1}(start = N \mid location) = \frac{C(start=N|location)}{C(location)+|start\ values|} = \frac{2+1}{6+6} = \frac{3}{12} = 0.25$$

$$- p_{+1}(start = E \mid location) = \frac{C(start=E|location)}{C(location)+|start\ values|} = \frac{1+1}{6+6} = \frac{2}{12} = 0.17$$

$$- p_{+1}(start = A \mid location) = \frac{C(start=A|location)}{C(location)+|start\ values|} = \frac{0+1}{6+6} = \frac{1}{12} = 0.08$$

$$- p_{+1}(start = Z \mid location) = \frac{C(start=Z|location)}{C(location)+|start\ values|} = \frac{0+1}{6+6} = \frac{1}{12} = 0.08$$

$$- \text{But } (3 + 1 + 2 + 3 + 2 + 1 + 1)/12 = 13/12!$$

– We had to add $|start\ values|$ to the denominator to keep the PD valid, but we failed.

– We should have accounted for Z by adding it to the set of start values

– So how many unseen start values are there? $A-Z$? Plus $a-z$? Plus $A-\Omega$?

– Sometimes we don’t know how many possible values there are. Like when its all of the words.

– Add- λ smoothing doesn’t handle this nicely.

– Solution? Use the set of known values and hope for the best.

6 Features occurring more than once per instance

- Up to this point we have only worked with instances that have exactly one value of each feature per instance.
- But sometimes we want to classify data that has feaures that appear multiple times per instance, possibly with the same values.
- A canonical example of this is when we are classifying text and we simply use the words in the text as features.

- The dataset:

the movie was great and the acting was great , so i loved it, positive
i loved the movie despite the terrible acting, positive
the plot was great , but overall it was terrible, negative
that movie was terrible, negative
the movie had a terrible ending, but was great anyway, positive
what a terrible movie, negative
i've never seen anything so terrible, negative

- Could have stopwords removed and be transformed to:

word=movie,word=great,word=acting,word=great,word=loved,5words=true,positive
word=loved,word=movie,word=terrible,word=acting,5words=true,positive
word=plot,word=great,word=terrible,5words=true,negative
word=movie,word=terrible,5words=false,negative
word=movie,word=terrible,word=ending,word=great,5words=true,positive
word=terrible,word=movie,5words=false,negative
word=terrible,5words=true,negative

- 2 features: *word* and *5words*

Counting with duplicate features

- If a feature can appear multiple times per instance, we have to count them all
- So the full counts are $C(\text{feature} = F, \text{label} = L)$.
- But we also need to make our probability calculation more precise
- Before, we said: $p(\text{feature} = F \mid \text{label} = L) = \frac{C(\text{feature}=F, \text{label}=L)}{C(\text{label}=L)}$
 - This would mean, for our data:
 - $p(\text{word} = \text{great} \mid \text{label} = \text{positive}) = \frac{C(\text{word}=\text{great}, \text{label}=\text{positive})}{C(\text{label}=\text{positive})} = \frac{3}{3} = 1.0$
 - This would mean that if the instance is *positive*, then it *definitely* has “great” in it... which is false.
 - Moreover, this would not make a valid probability distribution, since there are other words that also appear with *positive*.
 - The probability is wrong because we counted multiple *word* features per instance, but only one *label* feature per instance.
- Since there is no longer one feature of each time per instance, we have to make our probability calculations more precise.
- $p(\text{feature} = F \mid \text{label} = L) = \frac{C(\text{feature}=F, \text{label}=L)}{\sum_f C(\text{feature}=f, \text{label}=L)}$
 - So we are no longer dividing the number of instance with the label and feature by the number of instances with the label.

- Instead, we have to count the number of (label,value) pairs extracted from the data for the particular feature, including multiple identical pairs in the same instance
- The probability $p(\text{feature} = F \mid \text{label} = L)$ is the proportion of times the particular feature was F in an instance with label L out of the total number of times the particular feature appeared in an instance with label L
- $p(\text{word} = \text{great} \mid \text{label} = \text{positive}) = \frac{C(\text{word}=\text{great},\text{label}=\text{positive})}{\sum_f C(\text{word}=f,\text{label}=\text{positive})} = \frac{3}{13} = 0.23$
 $p(\text{word} = \text{great} \mid \text{label} = \text{negative}) = \frac{C(\text{word}=\text{great},\text{label}=\text{negative})}{\sum_f C(\text{word}=f,\text{label}=\text{negative})} = \frac{1}{8} = 0.13$
- Note that we still calculate the prior the same way: percentage of *instances* that have the particular label.

7 Evaluating

- How do we tell if our classifier works?
 - Evaluating on the training data doesn't tell us about performance on new data
- How do we set hyper-parameters (e.g. λ) without cheating?
 - Tuning them to the test set wouldn't tell us about performance on new data

Data Split

- Divide labeled data into three groups: **train**, **development (dev)**, and **test**
- Use only the *train* data to set parameters (take counts)
- Test against the *dev* data using various choices for hyper-parameters (e.g. λ) to figure out the best settings.
- When all development is done, and you are ready to get final results for publication, evaluate *just once* on the *test* data, and record the results as your final results.

Data Size

- Larger training set means that we will have better count estimates and, therefore, better model parameter estimates
- Larger dev and test sets mean that we can trust the results more

Accuracy

- Proportion of test instances that are labeled correctly: $\frac{C(\text{correct})}{C(\text{total})}$
- 100 instances, 70 correctly labeled: 70% correct
- Simple and intuitive, but can be misleading

- 100 instances, 80 are *not spam*
 - If the system just always outputs *not spam*, then accuracy is 80%
 - But the model is worthless!

Precision and Recall

- Precision: Out of those that we labeled “spam”, what proportion are actually spam?
 $p(\text{is spam} \mid \text{says “spam”})$
- Recall: Out of those that are actually spam, what proportion did we correctly label as “spam”?
 $p(\text{says “spam”} \mid \text{is “spam”})$
- Categorize kinds of correctness and mistakes for *spam*
 - true positive (tp): model says “spam”, actually spam
 - false positive (fp): model says “spam”, actually not spam
 - false negative (fn): model says “not spam”, actually spam
 - true negative (tn): model says “not spam”, actually not spam
- $\text{precision} = \frac{tp}{tp+fp}$ $\text{recall} = \frac{tp}{tp+fn}$
- F_1 -score: A way of combining precision and recall into one number (using the harmonic mean)
 - $F_1 = 2 \cdot \frac{P \cdot R}{P+R}$
- 100 instances:
 - P & R on *spam*

		model says	
		spam	not spam
truth	spam	TP = 20	FN = 10
	not spam	FP = 15	TN = 55

- * Accuracy is 0.75
- * $\text{Precision} = \frac{tp}{tp+fp} = \frac{20}{20+15} = 0.57$
- * $\text{Recall} = \frac{tp}{tp+fn} = \frac{20}{20+10} = 0.67$
- * $F_1 = 2 \cdot \frac{P \cdot R}{P+R} = 2 \cdot \frac{0.57 \cdot 0.67}{0.57+0.67} = 0.62$

- P & R on *not spam*

		model says	
		spam	not spam
truth	spam	TN = 20	FP = 10
	not spam	FN = 15	TP = 55

- * Accuracy is 0.75
- * $\text{Precision} = \frac{tp}{tp+fp} = \frac{55}{55+10} = 0.85$
- * $\text{Recall} = \frac{tp}{tp+fn} = \frac{55}{55+15} = 0.79$
- * $F_1 = 2 \cdot \frac{P \cdot R}{P+R} = 2 \cdot \frac{0.85 \cdot 0.79}{0.85+0.79} = 0.89$

- 100 instances, 80 are *not spam*, and the system just always outputs *not spam*

– P & R on *spam*

		model says	
		spam	not spam
truth	spam	TP = 0	FN = 20
	not spam	FP = 0	TN = 80

* Accuracy is 0.80

* Precision = $\frac{tp}{tp+fp} = \frac{0}{0+0} = NaN$

* Recall = $\frac{tp}{tp+fn} = \frac{0}{0+20} = 0.0$

* $F_1 = 2 \cdot \frac{P \cdot R}{P+R} = ?$

– P & R on *not spam*

		model says	
		spam	not spam
truth	spam	TN = 0	FP = 20
	not spam	FN = 0	TP = 80

* Accuracy is 0.80

* Precision = $\frac{tp}{tp+fp} = \frac{80}{80+20} = 0.80$

* Recall = $\frac{tp}{tp+fn} = \frac{80}{80+0} = 1.00$

* $F_1 = 2 \cdot \frac{P \cdot R}{P+R} = 2 \cdot \frac{0.80 \cdot 1.00}{0.80+1.00} = 0.89$

8 Semi-Supervised Learning for Naïve Bayes

- At this point, we know how to train a naïve Bayes classifier from fully-supervised data (every training instance has a label)
- Can we use (lots of) unlabeled data to improve a NB classifier trained on just a small amount of labeled data?
- Can we learn a NB classifier with out any labeled data?

An excellent tutorial that carefully explains everything necessary to implement Gibbs sampling for naïve Bayes:

- Philip Resnik and Eric Hardisty. “Gibbs Sampling for the Uninitiated”. June 2010.
<http://www.cs.umd.edu/~hardisty/papers/gsfu.pdf>

Partial Supervision: The Idea

- We know how to use NB to classify if we have labeled data
- If we have unlabeled data, we can assign labels to all the instances, then it because labeled data, and we can use it for NB

- We don't want to just assign labels and assume that they are correct
- So we will make initial guesses for those labels, and iteratively revise them
- Critically, we won't revise by simply picking the most likely label. We will choose probabilistically, based on the relative probabilities of the different labels.

The Algorithm

1. Make a NB classifier from any labeled data available.
 2. For each unlabeled instance:
 - a. Use this classifier to calculate the probabilities of each label.
 - b. Make a probability distribution out of all label probabilities (normalize them so that they add up to 1.0)
 - c. Sample from the probability distribution by choosing a label with likelihood according to its probability
 - d. Assign that chosen label as the label for the instance
 3. For some large number of iterations, for each unlabeled instance:
 - a. Make a NB classifier from all the instances and their current labels.
 - c. Use this classifier to calculate the probabilities of each label.
 - c. Make a probability distribution out of all label probabilities (normalize them so that they add up to 1.0)
 - d. Sample from the probability distribution by choosing a label with likelihood according to its probability
 - e. Assign that chosen label as the revised label for the instance
 - f. Copy the instance with its new label into a "repository" of labeled instances.
 4. Make a final NB classifier from all of the labeled instances in the "repository" that we built up while we were sampling.
- We don't revise the labels of the supervised instances because we assume that they are correct (unless we don't trust them, or think they are out-of-date).
 - Always making probabilistic label choices prevents us from just accepting the MLE and getting stuck.
 - The choices will allow us to "walk" around the space of labels, where most of the walking happens in area that lead to better parameters.

Why does this work?

- Basically, it learns how group associated feature values together, and how to push different labels apart.
- We start with a "best guess" based on information we have, and revise our guesses over time.

- We don't just walk "uphill" based on the MLE.
- Instead we allow ourselves to explore lots of label choices, some of which may be very different from the MLE.
- As we become more confident in a judgement, its label is less likely to change because more of the similar instances will be using that label, which gives it a high probability.

Markov Chain Monte Carlo

- Markov: Uses the "Markov principle" which says to make decisions based only on recent information
 - We use only the current set of labels, not the full history of labels.
- Chain: Start with an initial set of labels, use it to change one label, then move to the next label and change it based on the new set of labels. Repeat for a long time.
- Monte Carlo: Each step in the chain has chance involved
 - We choose the next label by chance according to the relative probabilities of the labels.
 - Named for Monte Carlo in Monaco and the Monte Carlo Casino

9 Clustering

- Until now, all our classifiers have been learned from fully-supervised data (every training instance has a label)
- Can we learn a classifier when we have no labels at all?
- **Clustering** is when we try to automatically induce groupings without labels.

K-Means

1. Plot all instances in space
2. Choose K of them to be initial centroids
3. For each instance, associate it with its nearest centroid
4. Update each centroid to the centroid of its currently-associated instances
5. Repeat step 3-4 until the centroids stop moving
6. For new instances, classify based on their nearest centroid