

# Ling 5801: Lecture Notes 13

## Probability Models

### Contents

13.1 Full joint models and sparse data problems . . . . .	1
13.2 Marginals and conditionals . . . . .	2
13.3 Factored models . . . . .	3
13.4 Graphical representation ('Bayes net') . . . . .	4
13.5 Conditional probability model class . . . . .	4
13.6 Estimating $\theta$ from fully-specified ('annotated') data . . . . .	7
13.7 Inducing $\theta$ from not-fully-specified ('unannotated') data . . . . .	7
13.8 Independence assumptions . . . . .	7
13.9 Another example: speech components . . . . .	8
13.10 Generative vs discriminative models . . . . .	9

### 13.1 Full joint models and sparse data problems

We can estimate probabilities for new atomic events based on frequencies of these events in training (assume a probability space  $\langle O, 2^O, P \rangle$  and frequency space  $\langle O, 2^O, F \rangle$  with  $O = X_1 \times X_2 \times \dots \times X_V$ ):

$$P(x_1, x_2, \dots, x_V) \stackrel{\text{def}}{=} \frac{F(x_1, x_2, \dots, x_V)}{F(O)}$$

This is called a **full joint distribution**.

For example, if we have a space of regions, words and sounds:  $O = \text{Reg} \times \text{Wrd} \times \text{Acou}$ , where:

- $\text{Reg} = \{\text{ohio, phil, ...}\}$  (speaker region)
- $\text{Wrd} = \{\text{neck, knack, ...}\}$  (speaker's intended word)
- $\text{Acou} = \{[\text{n}\epsilon\text{k}], [\text{n}\ae\text{k}], \dots\}$  (listener's observed phone)

For example, Ohians never pronounce /ε/ as [æ], but Philadelphians often do:

- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{ohio}, \text{Wrd} = \text{knack}, \text{Acou} = [\text{n}\ae\text{k}]) = .2$
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{ohio}, \text{Wrd} = \text{knack}, \text{Acou} = [\text{n}\epsilon\text{k}]) = 0$
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{ohio}, \text{Wrd} = \text{neck}, \text{Acou} = [\text{n}\ae\text{k}]) = 0$  (never)
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{ohio}, \text{Wrd} = \text{neck}, \text{Acou} = [\text{n}\epsilon\text{k}]) = .3$
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{phil}, \text{Wrd} = \text{knack}, \text{Acou} = [\text{n}\ae\text{k}]) = .2$
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{phil}, \text{Wrd} = \text{knack}, \text{Acou} = [\text{n}\epsilon\text{k}]) = 0$
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{phil}, \text{Wrd} = \text{neck}, \text{Acou} = [\text{n}\ae\text{k}]) = .2$  (often)
- $P_{\theta_{\text{pron}}}(\text{Reg} = \text{phil}, \text{Wrd} = \text{neck}, \text{Acou} = [\text{n}\epsilon\text{k}]) = .1$

We can write this as probability table, where rows sum to one:

$$P_{\theta_{pron}}(Reg, Wrd, Acou) =$$

ohio	ohio	ohio	ohio	phil	phil	phil	phil
knack	knack	neck	neck	knack	knack	neck	neck
[næk]	[nɛk]	[næk]	[nɛk]	[næk]	[nɛk]	[næk]	[nɛk]
.2	0	0	.3	.2	0	.2	.1

Each probability in the model is called a *parameter*.

There's one for each combination of values: in this case  $2 \cdot 2 \cdot 2 = 8$ .

But if there are too many outcomes for the training data, full joints can make arbitrary predictions.

For example, consider this (junk/not-junk) email classifier algorithm:

1. Memorize a set of training emails (random variables are for first word, second word, etc.).
2. Compare each test email to the training set (estimate probability of atomic event).
  - (a) If you find an exact match, report the class from the matching training example.
  - (b) If you find no exact match, call it whatever is most common in training data (not-junk).

This will very rarely find an exact match, and end up classifying virtually nothing as junk.

This is called a **sparse data problem**. It shows the importance of *generalizing* from data.

## 13.2 Marginals and conditionals

One way to make better use of scarce data is to consolidate training instances across conditions.

We can do this by marginalizing or summing out certain variables that don't matter as much.

1. We can define **marginal** distributions from the Kolmogorov axioms:

$$P(y, y', \dots) = \sum_{x, x', \dots \in X \times X' \times \dots} P(x, x', \dots, y, y', \dots)$$

For example:

$$P_{\theta_{pron}}(Wrd = w, Acou = a) = \sum_{r \in Reg} P_{\theta_{pron}}(Reg = r, Wrd = w, Acou = a)$$

$$P_{\theta_{pron}}(Wrd, Acou) =$$

knack	knack	neck	neck
[næk]	[nɛk]	[næk]	[nɛk]
.4	0	.2	.4

Another example:

$$P_{\theta_{pron}}(Acou = a) = \sum_{r \in Reg, w \in Wrd} P_{\theta_{pron}}(Reg = r, Wrd = w, Acou = a)$$

$$P_{\theta_{pron}}(Acou) = \begin{array}{|c|c|} \hline [næk] & [nɛk] \\ \hline .6 & .4 \\ \hline \end{array}$$

2. We can then define **conditional** distributions based on these marginals:

$$P(x, x', \dots | y, y', \dots) = \frac{P(x, x', \dots, y, y', \dots)}{P(y, y', \dots)}$$

For example:

$$\begin{aligned} P_{\theta_{pron}}(Wrd = w | Acou = a) &= \frac{P_{\theta_{pron}}(Wrd = w, Acou = a)}{P_{\theta_{pron}}(Acou = a)} \\ &= \frac{\sum_{r \in Reg} P_{\theta_{pron}}(Reg = r, Wrđ = w, Acou = a)}{\sum_{r \in Reg, w \in Wrđ} P_{\theta_{pron}}(Reg = r, Wrđ = w, Acou = a)} \end{aligned}$$

Table has one row for each combination of conditioned-on variable values:

$$P_{\theta_{pron}}(Wrd | Acou) = \begin{array}{|c|c|c|} \hline Acou & knack & neck \\ \hline [næk] & .667 & .333 \\ \hline [nɛk] & 0 & 1 \\ \hline \end{array}$$

So, by our model, we can calculate that [næk] is probably from knackin the general case:

$$P_{\theta_{pron}}(Wrd = knack | Acou = [næk]) = \frac{.2 + .2}{.2 + 0 + .2 + .2} = .667$$

but if the speaker is from Philadelphia, the odds are even:

$$P_{\theta_{pron}}(Wrd = knack | Acou = [næk], Reg = phil) = \frac{.2}{.2 + .2} = .5$$

## Practice

Write a Python program to calculate  $P(Wrd | Acou)$  in a dictionary  $WgivA[w, a]$  given a full joint distribution  $P(Reg, Wrđ, Acou)$  in a dictionary  $RWA[r, w, a]$ .

## 13.3 Factored models

We can consolidate data by first factoring our model, using the definition of conditional probability (you can validate this by cross-cancelling the numerators and denominators):

$$\begin{aligned} P(x_1, x_2, \dots, x_{V-1}, x_V) &= P(x_1) \cdot \underbrace{\frac{P(x_1, x_2)}{P(x_1)}}_{P(x_2 | x_1)} \cdot \underbrace{\frac{P(x_1, x_2, x_3)}{P(x_1, x_2)}}_{P(x_3 | x_1, x_2)} \dots \underbrace{\frac{P(x_1, x_2, \dots, x_{V-1}, x_V)}{P(x_1, x_2, \dots, x_{V-1})}}_{P(x_V | x_1, x_2, \dots, x_{V-1})} \\ &= P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_1, x_2) \dots P(x_V | x_1, x_2, \dots, x_{V-1}) \end{aligned}$$

This is called a **chain rule decomposition**.

Now, in our example:

$$\begin{aligned}
P_{\theta_{Pron}}(Reg, Wrd, Acou) &= \frac{P_{\theta_{Pron}}(Reg)}{1} \cdot \frac{P_{\theta_{Pron}}(Reg, Wrd)}{P_{\theta_{Pron}}(Reg)} \cdot \frac{P_{\theta_{Pron}}(Reg, Wrd, Acou)}{P_{\theta_{Pron}}(Reg, Wrd)} \\
&= P_{\theta_{Reg}}(Reg) \cdot P_{\theta_{Wrd}}(Wrd | Reg) \cdot P_{\theta_{Acou}}(Acou | Reg, Wrd)
\end{aligned}$$

Now the number of parameters in any variable's conditional probability distribution:

$$P_{\theta_{X_i}}(X_i | X_1, \dots, X_{i-1})$$

is the product of the cardinalities of its modeled and conditioned-on variables:

$$|\theta_{X_i}| = |X_i| \cdot |X_1| \cdot \dots \cdot |X_{i-1}|$$

In our example:

$$|\theta_{Acou}| = |Acou| \cdot |Reg| \cdot |Wrd| = 2 \cdot 2 \cdot 2 = 8 \text{ parameters}$$

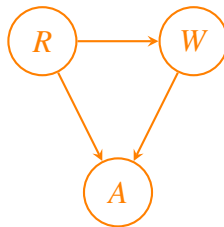
So far this is no better than the full joint distribution... but that's ok for now.

## 13.4 Graphical representation ('Bayes net')

Terms  $P_{\theta_{X_i}}(X_i | C_{X_i})$  from chain can be represented graphically —

- **circles** for random variables  $X_i$
- **arrows** for conditional dependencies (from conditioned-on  $C_{X_i}$  to modeled variables)

For example, probability space  $\langle R \times W \times A, 2^{R \times W \times A}, P \rangle$  would be drawn:



Graphical models are like sudoku: some values given, some unknown, values interdepend

In fact, they can be considered a *generalization* of sudoku puzzles!

## 13.5 Conditional probability model class

Extend 'model.py' to implement conditional model: dictionary of dictionary

(You may find this useful for certain problem set questions!)

```

import re
# define distribution to map value tuples to probs (or frequencies or scores)
class Model(dict):

```

```

# init with model id
def __init__(self,i=''):
    self.id = i
# read model
def read(self,s):
    m = re.search('^ *'+self.id+' +: +(.*?) += +(.*?) *',s)
    if m is not None:
        v = tuple(re.split(' +',m.group(1)))
        if len(v)==1: v = v[0]
        self[v] = float(m.group(2))
# write model
def write(self):
    for v in sorted(self):
        s = self.id
        s = s + ' :'
        if type(v) is tuple:
            for f in v:
                s = s + ' ' + f
        else: s = s + ' ' + v
        print( s + ' = ' + str(self[v]) )

# define model to map condition tuples to distributions
class CondModel(dict):
    # populate with default values when queried on missing keys
    def __missing__(self,k):
        self[k]=Model()
        return self[k]
    # define get without promiscuity
    def get(self,k):
        return dict.get(self,k,Model())
# init with model id
def __init__(self,i):
    self.id = i
# read model
def read(self,s):
    m = re.search('^ *'+self.id+' +(.*?) +: +(.*?) += +(.*?) *',s)
    if m is not None:
        c = tuple(re.split(' +',m.group(1)))
        if len(c)==1: c = c[0]
        v = tuple(re.split(' +',m.group(2)))
        if len(v)==1: v = v[0]
        self[c][v] = float(m.group(3))
# write model
def write(self):
    for c in sorted(self):
        for v in sorted(self[c]):
            s = self.id

```

```

    if type(c) is tuple:
        for f in c:
            s = s + ' ' + f
    else: s = s + ' ' + c
    s = s + ' :'
    if type(v) is tuple:
        for f in v:
            s = s + ' ' + f
    else: s = s + ' ' + v
    print( s + ' = ' + str(self[c][v]) )

```

Example, e.g. type into a program file 'myprog.py':

```

import re
import sys
import model

# read in params beginning with 'R' (prior model: no value before colon)
# read in params beginning with 'W' (prior model: no value before colon)
# read in params beginning with 'A' (cond model: value before colon)
R = model.Model('R')
W = model.Model('W')
A = model.CondModel('A')
for line in sys.stdin:
    R.read(line)
    W.read(line)
    A.read(line)

# use a prior model
for r in R:
    print ( 'prob of '+r+' is '+str(R[r]) )
# use a conditional model
for r,w in A:
    for a in A[r,w]:
        print ( 'prob of '+a+' given '+r+' and '+w+' is '+str(A[r,w][a]) )
# calc prob of W='knack' given R='phil' and A='[naek]'
probAnyW = 0.0
for w in W:
    probAnyW = probAnyW + R['phil'] * W[w] * A['phil',w]['[naek]']
print ( R['phil'] * W['knack'] * A['phil','knack']['[naek]'] / probAnyW )

```

Reads files in the following format:

```

R : ohio = .5
R : phil = .5
W : knack = .4
W : neck = .6
A ohio knack : [naek] = 1.0
A ohio neck : [nek] = 1.0

```

```

A phil knack : [naek] = 1.0
A phil neck  : [nek]  = .333333
A phil neck  : [naek] = .666667
:

```

then describes the information in the file and then prints the probability of knack:

```
0.499999875
```

### 13.6 Estimating $\theta$ from fully-specified ('annotated') data

Every  $\theta_{X_i}$  is a conditional probability table  $P_{\theta_{X_i}}(X_i | C_{X_i})$ .

Simply count instances of  $x_1, \dots, x_i \in C_{X_i} \times X_i$  and divide by count of  $x_1, \dots, x_{i-1} \in C_{X_i}$ .

This is called **relative frequency estimation**.

### 13.7 Inducing $\theta$ from not-fully-specified ('unannotated') data

You'll study this in comp ling 2!

### 13.8 Independence assumptions

If we then think some variables don't influence others much, we can remove them from conditions:

$$P_{\theta_X}(x|y, y', \dots, z, z', \dots) \stackrel{\text{def}}{=} P(x|y, y', \dots)$$

These re-definitions are called **independence assumptions**.

For example, our model of pronunciation variation may become:

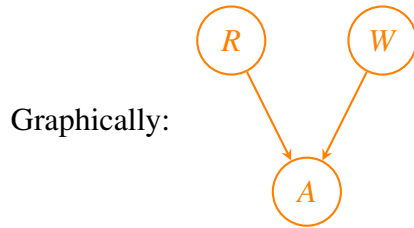
$$\begin{aligned}
&\langle \text{Reg} \times \text{Wrd} \times \text{Acou}, 2^{\text{Reg} \times \text{Wrd} \times \text{Acou}}, P \rangle \\
&P_{\theta_{\text{Reg}}}(\text{Reg}) \stackrel{\text{def}}{=} P(\text{Reg}) \\
&P_{\theta_{\text{Wrd}}}(\text{Wrd} | \text{Acou}) \stackrel{\text{def}}{=} P(\text{Wrd}) \\
&P_{\theta_{\text{Acou}}}(\text{Acou} | \text{Reg}, \text{Wrd}) \stackrel{\text{def}}{=} P(\text{Acou} | \text{Reg}, \text{Wrd})
\end{aligned}$$

The joint probability is:

$$P_{\theta_{\text{Pron}}}(\text{Reg}, \text{Wrd}, \text{Acou}) = P_{\theta_{\text{Reg}}}(\text{Reg}) \cdot P_{\theta_{\text{Wrd}}}(\text{Wrd} | \text{Reg}) \cdot P_{\theta_{\text{Acou}}}(\text{Acou} | \text{Reg}, \text{Wrd})$$

$$P_{\theta_{\text{Pron}}}(\text{Reg}, \text{Wrd}, \text{Acou}) \stackrel{\text{def}}{=} P_{\theta_{\text{Reg}}}(\text{Reg}) \cdot P_{\theta_{\text{Wrd}}}(\text{Wrd}) \cdot P_{\theta_{\text{Acou}}}(\text{Acou} | \text{Reg}, \text{Wrd})$$

because of independence assumption:  $P_{\theta_{Pron}}(Wrd | Reg) \stackrel{\text{def}}{=} P_{\theta_{Pron}}(Wrd)$



### 13.9 Another example: speech components

Here's an example of speech components: Phone, Voice, Back, Formant frequencies (binned):

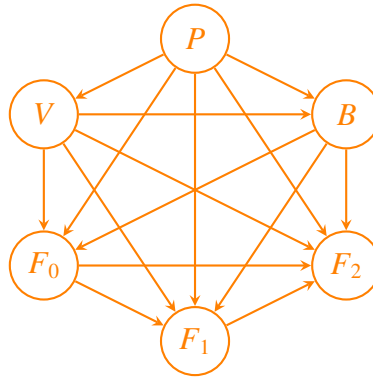
$$\langle P \times V \times B \times F_0 \times F_1 \times F_2, 2^{P \times V \times B \times F_0 \times F_1 \times F_2}, P \rangle$$

where  $P = \{/i/, /u/\}, V = \{+, -\}, B = \{+, -\}, F_0 = \mathbb{I}_0^{99}, F_1 = \mathbb{I}_0^{99}, F_2 = \mathbb{I}_0^{99}$

1. Chain rule decomposition (no independence assumptions):

$$P_{\theta_{Sp}}(P, V, B, F_0, F_1, F_2) = P_{\theta_P}(P) \cdot P_{\theta_V}(V | P) \cdot P_{\theta_B}(B | P, V) \cdot P_{\theta_{F_0}}(F_0 | P, V, B) \cdot P_{\theta_{F_1}}(F_1 | P, V, B, F_0) \cdot P_{\theta_{F_2}}(F_2 | P, V, B, F_0, F_1)$$

Graphed, it produces this evil hexagram:



Here  $|\theta_{F_2}| = 100 \cdot 2 \cdot 2 \cdot 2 \cdot 100 \cdot 100 = 8,000,000$  parameters!

2. With independence assumptions:

$$P_{\theta_P}(P) \stackrel{\text{def}}{=} P(P)$$

$$P_{\theta_V}(V | P) \stackrel{\text{def}}{=} P(V | P)$$

$$P_{\theta_B}(B | P, V) \stackrel{\text{def}}{=} P(B | P)$$

$$P_{\theta_{F_0}}(F_0 | P, V, B) \stackrel{\text{def}}{=} P(F_0 | V)$$

$$P_{\theta_{F_1}}(F_1 | P, V, B, F_0) \stackrel{\text{def}}{=} P(F_1 | B)$$

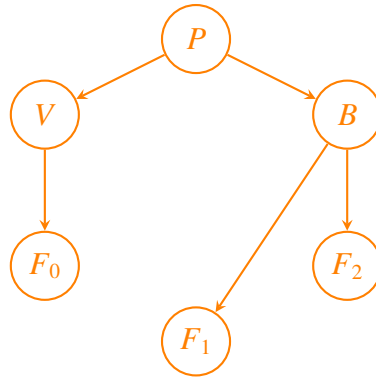
$$P_{\theta_{F_2}}(F_2 | P, V, B, F_0, F_1) \stackrel{\text{def}}{=} P(F_2 | B)$$



and joint distribution:

$$P_{\theta_{Sp}}(P, V, B, F_0, F_1, F_2) \stackrel{\text{def}}{=} P_{\theta_P}(P) \cdot P_{\theta_V}(V | P) \cdot P_{\theta_B}(B | P) \cdot P_{\theta_{F_0}}(F_0 | V) \cdot P_{\theta_{F_1}}(F_1 | B) \cdot P_{\theta_{F_2}}(F_2 | B)$$

it looks less busy:



Now  $|\theta_{F_2}| = 100 \cdot 2 = 200$  parameters!

### Practice

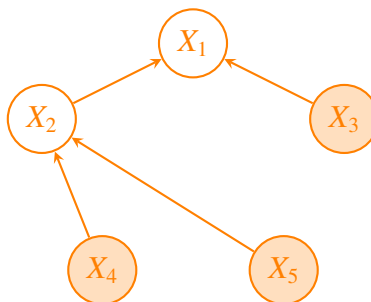
1. Draw a model of pet communication given random variables for vocalization, hunger, location, and tail being stepped on. The pet vocalizes when it's hungry, or when its tail is stepped on. Its tail gets stepped on only when it is in the kitchen.
2. Write the factored equation for  $P(\text{Hunger}, \text{TailStep} | \text{Voc}, \text{Loc})$
3. Write a program to do this calculation in `modHTgivVL[v, l][h, t]` given models `modL[l], modT[l][t], modH[h], modV[t, h][v]`

## 13.10 Generative vs discriminative models

Random variables that are observed at test time are typically shaded in graphical representations.

These are called **observations**; the others are called **hidden** variables.

When these observed variables are conditioned on, the model is called **discriminative**:



When observed variables are *not* conditioned on, the model is called **generative**:

