

Probability Basics, Density Estimation

These slides were assembled by Byron Boots, with only minor modifications from Eric Eaton's slides and grateful acknowledgement to the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution.

Robot Image Credit: Viktoriya Sukhanova © 123RF.com

The Joint Distribution

Recipe for making a joint distribution of *d* variables:

- Make a probability table listing all combinations of values of your variables (if there are *d* Boolean variables then the table will have 2^d rows).
- 1. For each combination of values, say how probable it is.
- If you subscribe to the axioms of probability, those numbers must sum to 1.

e.g., Boolean variables A, B, C

Α	В	С	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Slide © Andrew Moore

Inferring Marginal Probabilities from the Joint

	alarm		−alarm	
	earthquake	−earthquake	earthquake	¬earthquake
burglary	0.01	0.08	0.001	0.009
¬burglary	0.01	0.09	0.01	0.79

 $P(alarm) = \sum_{b,e} P(alarm \land \text{Burglary} = b \land \text{Earthquake} = e)$ = 0.01 + 0.08 + 0.01 + 0.09 = 0.19

 $P(burglary) = \sum_{a,e} P(\text{Alarm} = a \land burglary \land \text{Earthquake} = e)$ = 0.01 + 0.08 + 0.001 + 0.009 = 0.1

Conditional Probability

• $P(A \mid B) = Probability that A is true given B is true$



What if we already know that *B* is true?

That knowledge changes the probability of *A*

Because we know we're in a world where *B* is true

$$P(A \mid B) = \frac{P(A \land B)}{P(B)}$$
$$P(A \land B) = P(A \mid B) \times P(B)$$

Example: Conditional Probabilities

$$P(A \mid B) = \frac{P(A \land B)}{P(B)}$$
$$P(A \land B) = P(A \mid B) \times P(B)$$

	alarm	−alarm
burglary	0.09	0.01
¬burglary	0.1	0.8

P(Alarm, Burglary) =

P(burglary alarm)	= P(burglary \land alarm) / P(alarm)
	= 0.09 / 0.19 = 0.47

P(alarm | burglary) = P(burglary \land alarm) / P(burglary) = 0.09 / 0.1 = 0.9

P(burglary \land alarm) = P(burglary | alarm) P(alarm) = 0.47 * 0.19 = 0.09

Example: Inference from Conditional Probability

$$P(A \mid B) = \frac{P(A \land B)}{P(B)}$$
$$P(A \land B) = P(A \mid B) \times P(B)$$



P(headache) = 1/10 P(flu) = 1/40P(headache | flu) = 1/2

"Headaches are rare and flu is rarer, but if you're coming down with the flu, then there's a 50-50 chance you'll have a headache."

Example: Inference from Conditional Probability

$$P(A \mid B) = \frac{P(A \land B)}{P(B)}$$
$$P(A \land B) = P(A \mid B) \times P(B)$$



P(headache) = 1/10 P(flu) = 1/40P(headache | flu) = 1/2

One day you wake up with a headache. You think: "Drat! 50% of flus are associated with headaches so I must have a 50-50 chance of coming down with flu."

Is this reasoning good?

Example: Inference from Conditional Probability

$$P(A \mid B) = \frac{P(A \land B)}{P(B)}$$
$$P(A \land B) = P(A \mid B) \times P(B)$$

$$P(headache) = 1/10$$
Want to solve for: $P(flu) = 1/40$ $P(headache \land flu) = ?$ $P(headache \mid flu) = 1/2$ $P(flu \mid headache) = ?$

P(headache
$$\land$$
 flu) = P(headache | flu) x P(flu)
= 1/2 x 1/40 = 0.0125

- P(flu | headache)
- = $P(headache \land flu) / P(headache)$ = 0.0125 / 0.1 = 0.125

Based on example by Andrew Moore

Bayes' Rule

$$P(A \mid B) = \frac{P(B \mid A) \times P(A)}{P(B)}$$

- Exactly the process we just used
- The most important formula in probabilistic machine learning

(Super Easy) Derivation: $P(A \land B) = P(A \mid B) \times P(B)$ $P(B \land A) = P(B \mid A) \times P(A)$ these are the same Just set equal...

$$P(A \mid B) \times P(B) = P(B \mid A) \times P(A)$$
 and solve...



Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418**

Bayes' Rule

- Allows us to reason from evidence to hypotheses
- Another way of thinking about Bayes' rule:

$$P(\text{hypothesis} \mid \text{evidence}) = \frac{P(\text{evidence} \mid \text{hypothesis}) \times P(\text{hypothesis})}{P(\text{evidence})}$$

In the flu example:

P(headache) = 1/10 P(flu) = 1/40

P(headache | flu) = 1/2

Given evidence of headache, what is P(flu | headache) ?

Solve via Bayes rule!

Independence

- When two sets of propositions do not affect each others' probabilities, we call them **independent**
- Formal definition:

$$A \bot B \quad \leftrightarrow \quad P(A \land B) = P(A) \times P(B)$$
$$\quad \leftrightarrow \quad P(A \mid B) = P(A)$$

For example, {moon-phase, light-level} might be independent of {burglary, alarm, earthquake}

- Then again, maybe not: Burglars might be more likely to burglarize houses when there's a new moon (and hence little light)
- But if we know the light level, the moon phase doesn't affect whether we are burglarized

Exercise: Independence

D(amout & study & puop)	sn	nart	smart		
P(smart ~ study ~ prep)	study	study	study ¬study		
prepared	0.432	0.16	0.084	0.008	
-prepared	0.048	0.16	0.036	0.072	

Is *smart* independent of *study*?

Is *prepared* independent of *study*?

Exercise: Independence

D(amout a study a prop)	sn	nart	smart		
P(smart ~ study ~ prep)	study	-study	study	study	
prepared	0.432	0.16	0.084	0.008	
-prepared	0.048	0.16	0.036	0.072	

Is smart independent of study? $P(study \land smart) = 0.432 + 0.048 = 0.48$ P(study) = 0.432 + 0.048 + 0.084 + 0.036 = 0.6 P(smart) = 0.432 + 0.048 + 0.16 + 0.16 = 0.8 $P(study) \times P(smart) = 0.6 \times 0.8 = 0.48$ Is prepared independent of study?

Conditional Independence

• Absolute independence of *A* and *B*:

$$A \bot B \quad \leftrightarrow \quad P(A \land B) = P(A) \times P(B)$$
$$\leftrightarrow \quad P(A \mid B) = P(A)$$

Conditional independence of *A* and *B* given *C* $A \perp B \mid C \iff P(A \land B \mid C) = P(A \mid C) \times P(B \mid C)$

- e.g., Moon-Phase and Burglary are *conditionally independent given* Light-Level
- This lets us decompose the joint distribution: $P(A \land B \land C) = P(A \mid C) \times P(B \mid C) \times P(C)$
 - Conditional independence is weaker than absolute independence, but still useful in decomposing the full joint

Take Home Exercise: Conditional independence

P(smart ^ study ^ prep) prepared ¬prepared	sn	nart	smart		
r(smart ~ study ~ prep)	study	study	study	study	
prepared	0.432	0.16	0.084	0.008	
-prepared	0.048	0.16	0.036	0.072	

Is *smart* conditionally independent of *prepared*, given *study*?

Is study conditionally independent of prepared, given smart?

Summary: Essential Probability Concepts • Marginalization: $P(B) = \sum_{v \in \text{values}(A)} P(B \land A = v)$

• Conditional Probability: $P(A \mid B) = \frac{P(A \land B)}{P(B)}$

• Bayes' Rule:
$$P(A \mid B) = \frac{P(B \mid A) \times P(A)}{P(B)}$$

Independence:

Density Estimation

How Can We Obtain a Joint Distribution?

Option 1: Elicit it from an expert human

- **Option 2:** Build it up from simpler probabilistic facts
- e.g, if we knew

P(a) = 0.7 P(b|a) = 0.2 P(b| \neg a) = 0.1 then, we could compute P(a \land b)

Option 3: Learn it from data...

Learning a Joint Distribution

Step 1:

Build a JD table for your attributes in which the probabilities are unspecified

Α	В	С	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

Step 2:

Then, fill in each row with:

 $\hat{P}(\text{row}) = \frac{\text{records matching row}}{\text{total number of records}}$

A	В	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Fraction of all records in which

A and B are true but C is false

Density Estimation

- Our joint distribution learner is an example of something called **Density Estimation**
- A Density Estimator learns a mapping from a set of attributes to a probability



Density Estimation

Compare it against the two other major kinds of models:



Slide © Andrew Moore

Evaluating Density Estimation

Test-set criterion for estimating performance on future data



Slide © Andrew Moore

Evaluating a Density Estimator

• Given a record **x**, a density estimator *M* can tell you how likely the record is:

 $\hat{P}(\mathbf{x} \mid M)$

- The density estimator can also tell you how likely the dataset is:
 - Under the assumption that all records were independently generated from the Density Estimator's JD (that is, i.i.d.)

$$\hat{P}(\mathbf{x}_1 \wedge \mathbf{x}_2 \wedge \ldots \wedge \mathbf{x}_n \mid M) = \prod_{i=1}^n \hat{P}(\mathbf{x}_i \mid M)$$
dataset

Example Small Dataset: Miles Per Gallon

From the UCI repository (thanks to Ross Quinlan)

• 192 records in the training set

mpg	modelyear	maker
good	75to78	asia
bad	70to74	america
bad	75to78	europe
bad	70to74	america
bad	70to74	america
bad	70to74	asia
bad	70to74	asia
bad	75to78	america
:	•	:
:	:	:
:	:	:
bad	70to74	america
good	79to83	america
bad	75to78	america
good	79to83	america
bad	75to78	america
good	79to83	america
good	79to83	america
bad	70to74	america
good	75to78	europe
bad	75to78	europe



Slide by Andrew Moore

Example Small Dataset: Miles Per Gallon

From the UCI repository (thanks to Ross Quinlan)

• 192 records in the training set

Slide by

				mpg	modelyear	maker			
mpg	modelyea	maker		bad	70to74	america	0.27551		
mpg	moderyea	IIIdKEI				asia	0.0255102		
good	75to78	asia				0.0000	0.0462064	-	
bad	70to74	america				europe	0.0155001		
$\hat{P}(\mathrm{da}$	tase	t <i>M</i>	$) = \prod_{i=1}^{n}$ = 3.	Ι 1 4	$\hat{P}(\mathbf{x}_i \times 10)$	_ <i>I</i> _20	Л) 93	(in this ca	use)
bad	75to78	america			13(0)1		0.0000122	-	
good	79t083	america				asia	0.0408163		
bad	70to74	america				europe	0.0357143		
aood	75to78	europe			78to83	america	0.112245		
bad	75to78	europe				asia	0.0714286		
Andrew Moore						europe	0.0357143		

Log Probabilities

For decent sized data sets, this product will underflow

$$\hat{P}(\text{dataset} \mid M) = \prod_{i=1}^{n} \hat{P}(\mathbf{x}_i \mid M)$$

• Therefore, since probabilities of datasets get so small, we usually use log probabilities

$$\log \hat{P}(\text{dataset} \mid M) = \log \prod_{i=1}^{n} \hat{P}(\mathbf{x}_i \mid M) = \sum_{i=1}^{n} \log \hat{P}(\mathbf{x}_i \mid M)$$

Example Small Dataset: Miles Per Gallon

From the UCI repository (thanks to Ross Quinlan)

• 192 records in the training set

					mpg	modelyear	maker			
	mpg	modelyear	maker]	bad	70to74	america	0.27551		
	good bad	75to78 70to74	asia america				asia europe	0.0255102		
log	${ m g}\hat{P}$	(dat	aset	<i>M</i>) =		$\sum_{i=1}^{n} 1$ -460	$\log 1$	$\hat{P}(\mathbf{x}_i \mid M)$ \hat{P} (in the set of th	1) his cas	e)
	bad	75to78 79to83	america america			75to77	america asia	0.0306122		
	good	79to83	america				eurone	0.0357143		
	bad	70to74	america]			curope	0.0007140		
	good	75to78	europe			78to83	america	0.112245		
	bad	75to78	europe				acia	0.0714286		
							asia	0.0114200		

Pros/Cons of the Joint Density Estimator

The Good News:

- We can learn a Density Estimator from data.
- Density estimators can do many good things...
 - Can sort the records by probability, and thus spot weird records (anomaly detection)
 - Can do inference
 - Ingredient for Bayes Classifiers (coming very soon...)

The Bad News:

• Density estimation by directly learning the joint is impractical, may result in adverse behavior

Curse of Dimensionality



Slide by Christopher Bishop

The Joint Density Estimator on a Test Set

	Set Size	Log likelihood
Training Set	196	-466.1905
Test Set	196	-614.6157

- An independent test set with 196 cars has a much worse log-likelihood
 - Actually it's a billion quintillion quintillion quintillion quintillion quintillion times less likely
- Density estimators can overfit...

...and the full joint density estimator is the overfittiest of them all!

Overfitting Density Estimators

