

Probability and Naive Bayes

Machine Learning
CS4824/ECE4424
Bert Huang
Virginia Tech

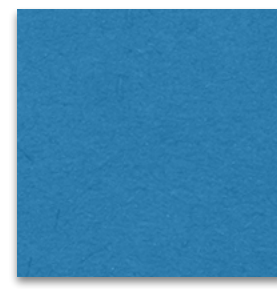
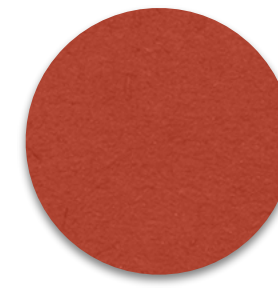
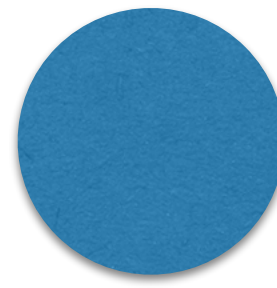
Outline

- Probabilistic identities
- Independence and conditional independence
- Naive Bayes
- Log tricks

Probability Identities

- Random variables in caps (**A**)
 - values in lowercase: **A = a** or just **a** for shorthand
- $P(a | b) = P(a, b) / P(b)$ conditional probability
- $P(a, b) = P(a | b) P(b)$ joint probability
- $P(b | a) = P(a | b) P(b) / P(a)$

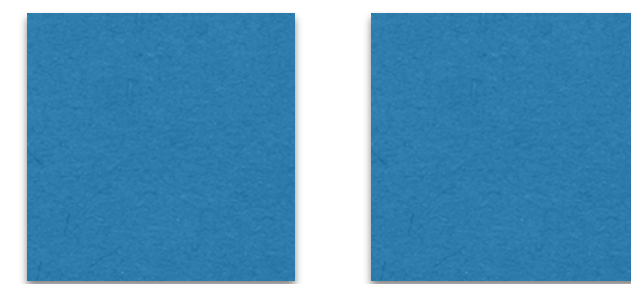
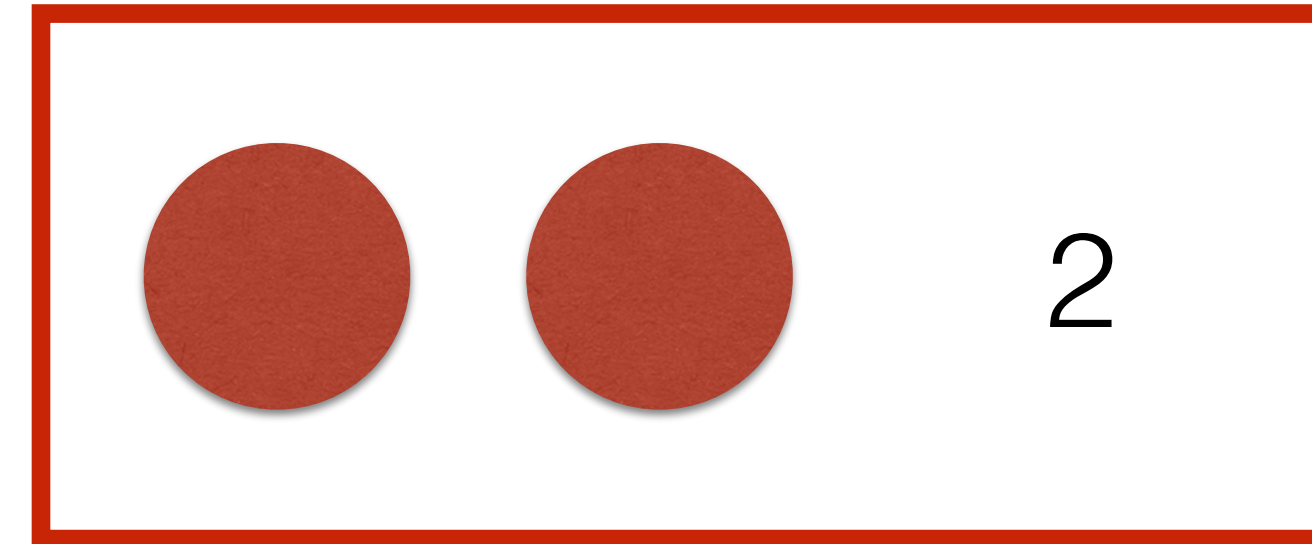
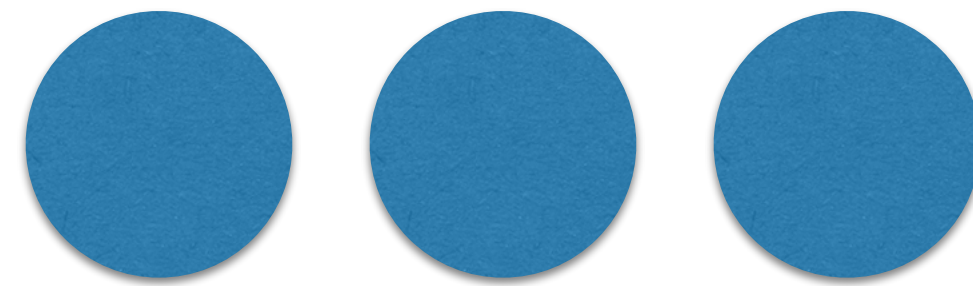
Probability via Counting



Probability via Counting

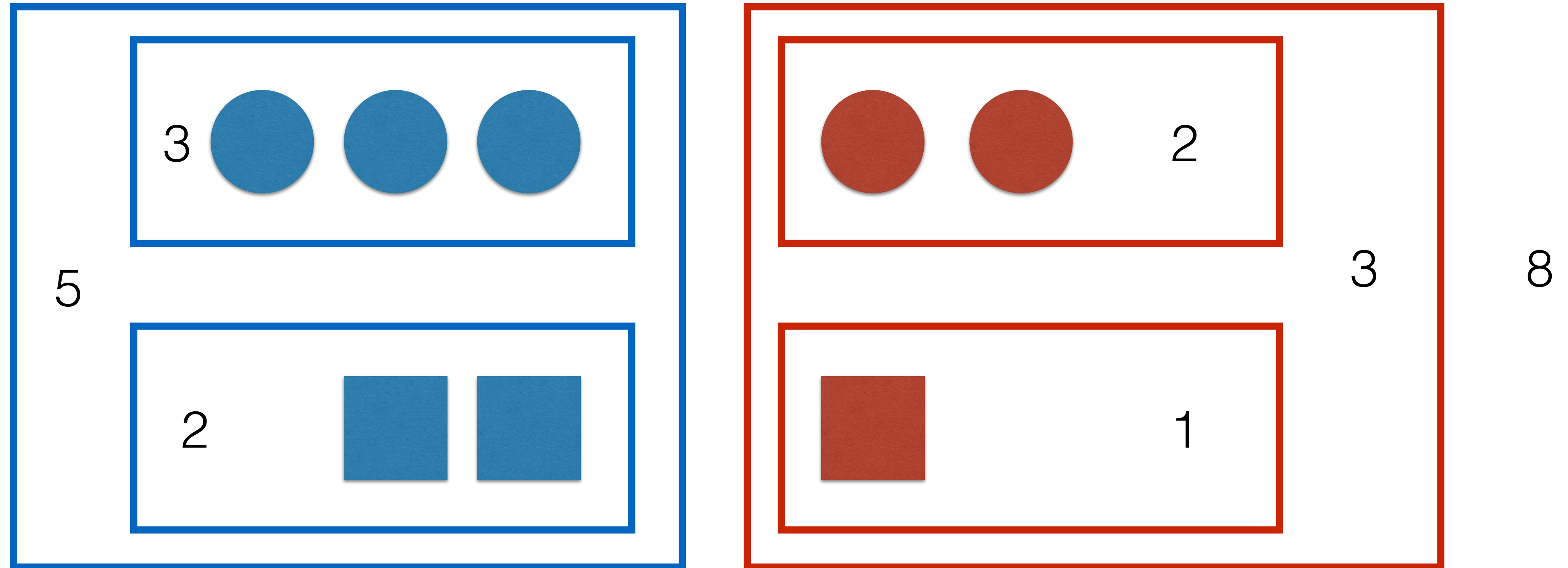
$P(\text{circle, red})$

$$= 2/8 = 0.25$$



8

Probability via Counting



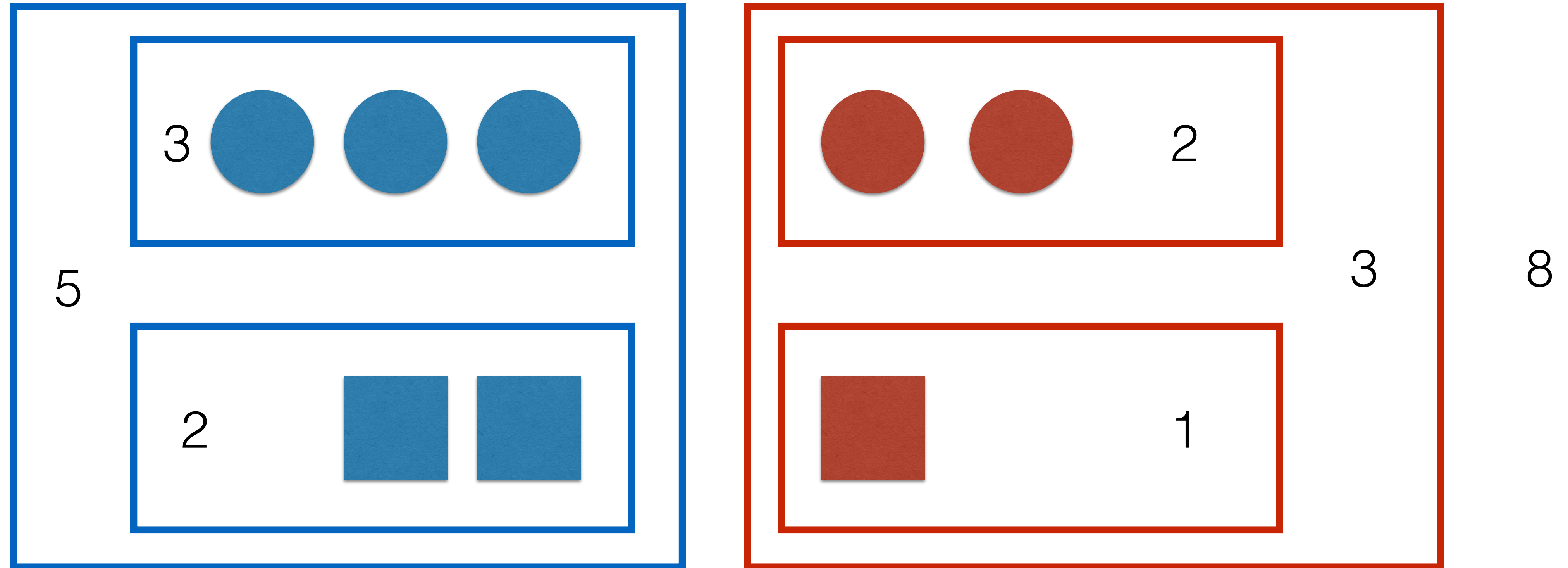
$$P(\text{circle} \mid \text{red}) = P(\text{circle, red}) / P(\text{red})$$

$2/3$

$2/8$

$3/8$

Probability via Counting



$$P(\text{circle} \mid \text{red}) P(\text{red}) = P(\text{circle, red})$$

$2/3$

$3/8$

$2/8$

Probability Identities

- Random variables in caps (**A**)
 - values in lowercase: **A = a** or just **a** for shorthand
- $P(a | b) = P(a, b) / P(b)$
- $P(a, b) = P(a | b) P(b)$
- $P(b | a) = P(a | b) P(b) / P(a)$

Bayes Rule

- $P(b | a)$
- $P(b | a) = P(a, b) / P(a)$
- $P(b | a) = P(a | b) P(b) / P(a)$

Classification

- $x \in \{0,1\}^d, y \in \{0,1\}$
- $f(x) \in \{0,1\}$
- Accuracy: $E[f(x) = y]$
- **Bayes optimal classifier:** $f(x) = \arg \max_y p(y | x)$
 - Seems natural, but why is this optimal?

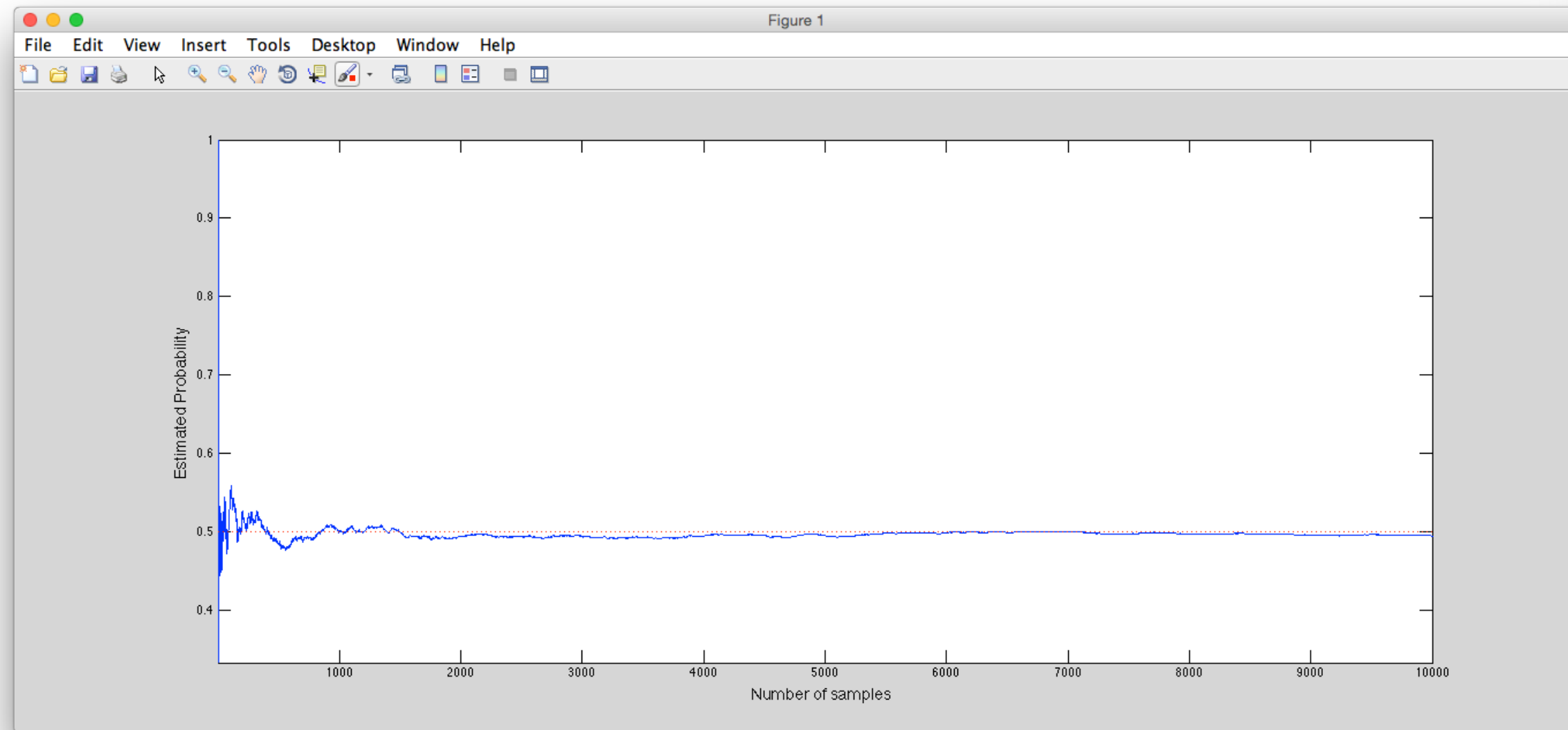
Back-of-Envelope for Bayes Optimal

- For each unique \mathbf{x} , $\mathbf{p}(\mathbf{y} \mid \mathbf{x})$ is a coin flip
- Assume we need \mathbf{n} samples to accurately estimate a coin flip
- How many unique \mathbf{x} 's? $n = 100$
 - $|\{0,1\}^d| = 2^d$ $d = 100$
- Need $\mathbf{n}2^d$ samples Need 1.2676506×10^{32} samples

1.2676506 x 10,000,000,000,000,000,000,000,000,000,000,000,000

How Many Samples?

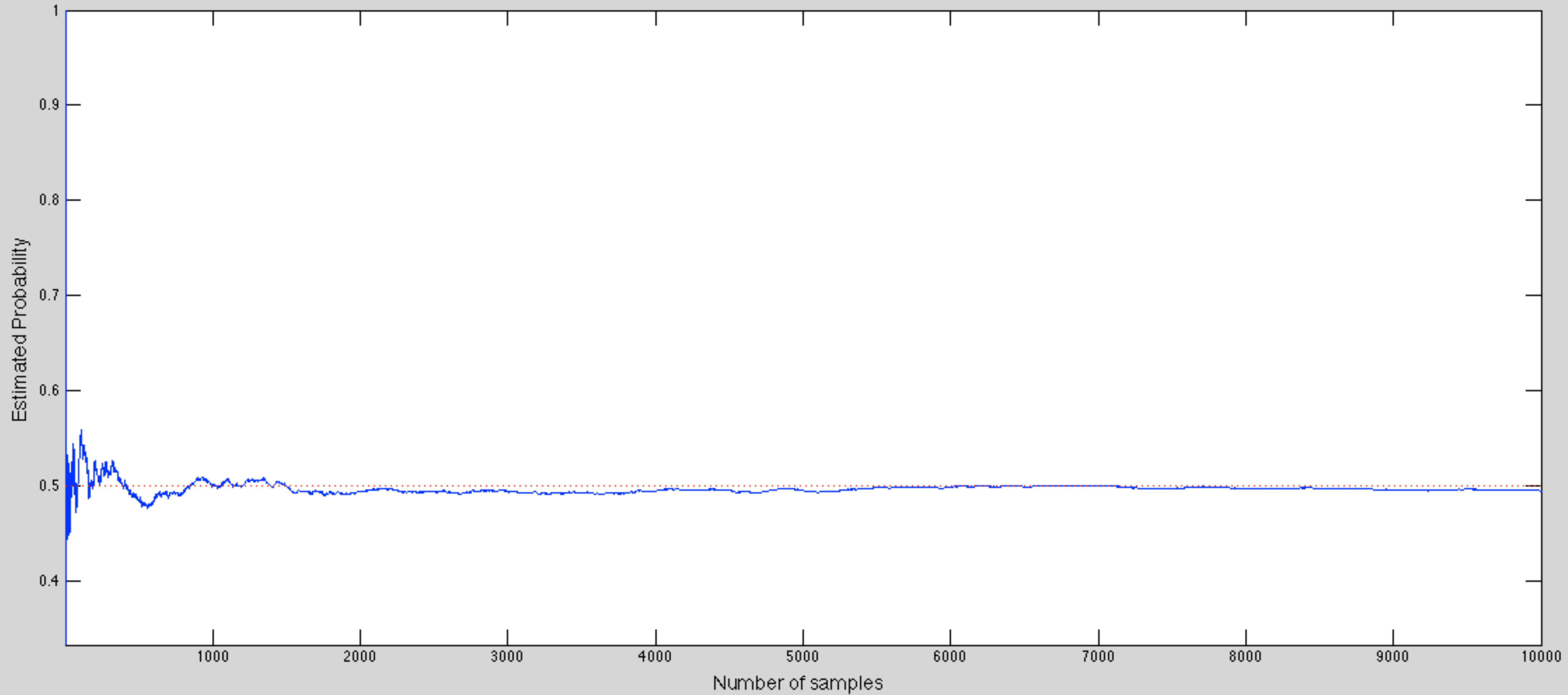
- Concentration bounds



```
Editor - /Users/bert/Dropbox/Teaching/ML/coinFlips.m
EDITOR PUBLISH VIEW
+ New + Open + Save + Find Files + Compare + Print + Insert + Comment + Indent + fx + % + % + % + Go To + Find + Breakpoints + Run + Run and Time + Run and Advance + Advance
FILE EDIT NAVIGATE BREAKPOINTS RUN
decisionTreePredict.m calculateInformationGain.m coinFlips.m
1 %% generate N coin flips
2
3 N = 10000;
4
5 heads = rand(N,1) > 0.5;
6
7 estimate = zeros(N,1);
8
9 for i = 1:N
10     estimate(i) = sum(heads(1:i)) / i;
11 end
12
13 plot(estimate);
14 xlabel('Number of samples', 'FontSize', 14);
15 ylabel('Estimated Probability', 'FontSize', 14);
16 hold on;
17 plot([1 N], [0.5 0.5], ':r');
18 hold off;
19 axis tight;
20
script Ln 19 Col 12
```

Figure 1

File Edit View Insert Tools Desktop Window Help



Independence

- **A** and **B** are independent iff $\mathbf{p(A, B) = p(A) p(B)}$
- **A** and **B** are *conditionally* independent given **C** iff $\mathbf{p(A, B | C) = p(A | C) p(B | C)}$

Naive Bayes

- Assume dimensions of \mathbf{x} are conditionally independent given \mathbf{y}
 - Bag of words: $p(\text{"virginia"}, \text{"tech"} \mid y) = p(\text{"virginia"} \mid y) p(\text{"tech"} \mid y)$
- $f(x) = \arg \max_y p(y \mid x)$
- $= \arg \max_y p(x \mid y) p(y) / p(x)$
- $= \arg \max_y p(x \mid y) p(y)$
- $= \arg \max_y p(y) \prod_j p(x_j \mid y)$

Bernoulli Maximum Likelihood

- $p(y) \prod_j p(x_j | y)$
- $p(Y = y) \leftarrow (\# \text{ examples where } Y = y) / (\# \text{ examples})$
- $p(X_j = x_j | y) \leftarrow (\# \text{ ex. where } Y = y \text{ and } X_j = x_j) / (\# \text{ ex. where } Y=y)$
- Learning by counting!

Breaking Maximum Likelihood

- Happy: "Great!"
- Happy: "Had a great day"
- Sad: ":-(Bad day"
- Sad: ":-(

	great	had	a	bad	day	:-(
Happy: "Great!"	1	0	0	0	0	0
Happy: "Had a great day"	1	1	1	0	1	0
Sad: ":-(Bad day"	0	0	0	1	1	1
Sad: ":-(0	0	0	0	0	1

- ???: "Had a bad day :-(

	great	had	a	bad	day	:-(
Happy	1.0	0.5	0.5	0.0	0.5	0
Sad	0.0	0.0	0.0	0.5	0.5	1.0

- $p(y) = 0.5$

- $p(\text{happy} \mid \dots) \propto 0.5 \times 0.0 \dots$

$$p(\text{sad} \mid \dots) \propto 0.5 \times 1.0 \times 0.0 \times \dots$$

Breaking Maximum Likelihood

- Happy: "Great!"
- Happy: "Had a great day"
- Sad: ":-(Bad day"
- Sad: ":-(

	great	had	a	bad	day	:-(
Happy: "Great!"	1	0	0	0	0	0
Happy: "Had a great day"	1	1	1	0	1	0
Sad: ":-(Bad day"	0	0	0	1	1	1
Sad: ":-(0	0	0	0	0	1

- ???: "Had a bad day :-(

	great	had	a	bad	day	:-(
Happy	1.0	0.5	0.5	0.0	0.5	0
Sad	0.0	0.0	0.0	0.5	0.5	1.0

- $p(y) = 0.5$

- $p(\text{happy} \mid \dots) \propto 0.5 \times 0.0 \dots$

$$p(\text{sad} \mid \dots) \propto 0.5 \times 1.0 \times 0.0 \times \dots$$

Fixing Maximum Likelihood

- $p(Z = z) \leftarrow (\# \text{ examples where } Z = z + \alpha) / (\# \text{ examples} + 2\alpha)$
 - E.g., $\alpha = 1$
- α vanishes as # of examples grows toward infinity
- When # is small, α prevents 1.0 or 0.0 estimates

Breaking Maximum Likelihood

- Happy: "Great!"
- Happy: "Had a great day"
- Sad: ":-(Bad day"
- Sad: ":-("

$$\frac{2 + 1}{2 + 2(1)} = 3/4$$

	great	had	a	bad	day	:-("
Happy: "Great!"	1	0	0	0	0	0
Happy: "Had a great day"	1	1	1	0	1	0
Sad: ":-(Bad day"	0	0	0	1	1	1
Sad: ":-("	0	0	0	0	0	1

- ???: "Had a bad day :-("

	great	had	a	bad	day	:-("
Happy	0.75	0.5	0.5	0.25	0.5	0.25
Sad	0.25	0.25	0.25	0.5	0.5	0.75

- $p(y) = 0.5$

- $p(\text{happy} \mid \dots) \propto 0.5 \times 0.25 \times 0.5 \times 0.5 \times 0.25 \times 0.5 \times 0.75 = 0.0029$
- $p(\text{sad} \mid \dots) \propto 0.5 \times 0.75 \times 0.25 \times 0.25 \times 0.5 \times 0.5 \times 0.75 = 0.0044$

Maximum a Posteriori

- Bernoulli: $p(Z | \theta) = \theta^Z (1 - \theta)^{(1-Z)}$
- Maximum likelihood: $\theta \leftarrow \arg \max_{\theta'} p(Z | \theta')$
- Maximum a posteriori = maximize posterior: $\theta \leftarrow \arg \max_{\theta'} p(\theta' | Z)$
- $p(\theta' | Z) = p(Z | \theta') p(\theta') / p(Z)$
- MAP: $\theta \leftarrow \arg \max_{\theta'} p(Z | \theta') p(\theta')$
- Previous trick equiv. to setting $p(\theta')$ to a Beta distribution

Continuous Data

- Conditional feature independence with continuous data?
- E.g., use normal distribution for $p(x_j | y)$
 - $p(y)$ is the same as before
 - $p(x_j | y)$ is the MLE for univariate normal

Log Tricks

- Each $p(x_j | y)$ is in $[0, 1]$
- Multiplying **d** of them quickly goes to numerical zero
 - E.g., $0.9^{256} = 1.932334983\text{E-}12$
- Instead, use log probabilities: $\log \prod_j p(x_j | y) = \sum_j \log p(x_j | y)$
 - E.g., $\log 0.9^{256} = 256 \log 0.9 = -11.71$

Summary

- Probabilistic identities
- Independence and conditional independence
- Naive Bayes
- Log tricks