# Machine Learning

Adapted from
Tom M. Mitchell
Carnegie Mellon University

Today:
- Bayes Rule
- Estimating parameters
  - MLE
  - MAP

some of these slides are derived
from William Cohen, Andrew
Moore, Aarti Singh, Eric Xing,
Carlos Guestrin.  - Thanks!

# Discriminative vs Generative Models

Goal is same

Calculate posterior distribution

Probability of Data, $\boldsymbol{P(y|X)}$

| Discriminative Models | Generative Models |
| --- | --- |
| Find the decision boundary that separates the classes | Say, you have two classes – $y_1$ and $y_2$, with features $X_1$ and $X_2$ |
| Only knows the differences between classes | First, looking at examples of $y_1$, build a model of what $y_1$ looks like/ distribution of $y_1$'s features |
| To classify a new example, see which side of the decision boundary it falls | Then, looking at examples of $y_2$, build a model of what $y_2$ looks like/ distribution of $y_2$'s features |
| | To classify a new example, match the new example with model of each class, to see whether the new example looks more like $y_1$ or more like $y_2$ we had seen in the training set |

# Discriminative vs Generative Models

Goal is same

Calculate posterior distribution

$$\hat{y} = \underset{x}{\text{argmax}} \, P(y|X)$$

Probability of Data, $P(y|X)$

| Discriminative Models | Generative Models |
|---|---|
| Find the decision boundary that separates the classes | Say, you have two classes – $y_1$ and $y_2$, with features $X_1$ and $X_2$ |
| Only knows the differences between classes | First, looking at examples of $y_1$, build a model of what $y_1$ looks like/ distribution of $y_1$'s features |
| To classify a new example, see which side of the decision boundary it falls | Then, looking at examples of $y_2$, build a model of what $y_2$ looks like/ distribution of $y_2$'s features |
| | To classify a new example, match the new example with model of each class, to see whether the new example looks more like $y_1$ or more like $y_2$ we had seen in the training set |

# Discriminative vs Generative Models

| Discriminative Models | Generative Models |
|---|---|
| Directly learn the function mapping $$h: X \rightarrow y$$ or, Calculate posterior distribution $$P(y\|X)$$ | Calculate posterior distribution $$P(y\|X)$$ HOW? |
| 1. Assume some functional form for $P(y\|X)$<br>2. Estimate parameters of $P(y\|X)$ directly from training data | |

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$ Bayes' rule
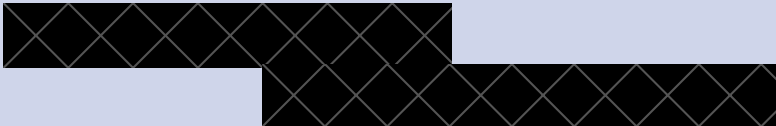


we call P(A) the "prior"

and P(A|B) the "posterior"

**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**

…by no means merely a curious speculation in the doctrine of chances, but necessary to be solved in order to a sure foundation for all our reasonings concerning past facts, and what is likely to be hereafter…. necessary to be considered by any that would give a clear account of the strength of *analogical* or *inductive reasoning…*

# Discriminative vs Generative Models

| Discriminative Models | Generative Models |
|---|---|
| Directly learn the function mapping $$h: X \rightarrow y$$ **or,** Calculate posterior distribution $$P(y\|X)$$ | Calculate posterior distribution $$P(y\|X)$$ from $P(X\|y)$ and $P(y)$ |
| 1. Assume some functional form for $P(y\|X)$ 2. Estimate parameters of $P(y\|X)$ directly from training data | |

# Other Forms of Bayes Rule

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$

$$P(A|B) = \frac{P(B \mid A)P(A)}{P(B \mid A)P(A) + P(B \mid \sim A)P(\sim A)}$$

$$P(A|B \wedge X) = \frac{P(B \mid A \wedge X)P(A \wedge X)}{P(B \wedge X)}$$

# Applying Bayes Rule

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

A = you have the flu,   B = you just coughed

Assume:
P(A) = 0.05
P(B|A) = 0.80
P(B| ~A) = 0.20

what is P(flu | cough) = P(A|B)?

$$P(A \mid B) = \frac{.8 \cdot .05}{.8 \cdot .05 + 0.2 \cdot 0.95} = 0.17$$

$$P(A) = 1 - P(\neg A)$$

what does all this have to do with function approximation?

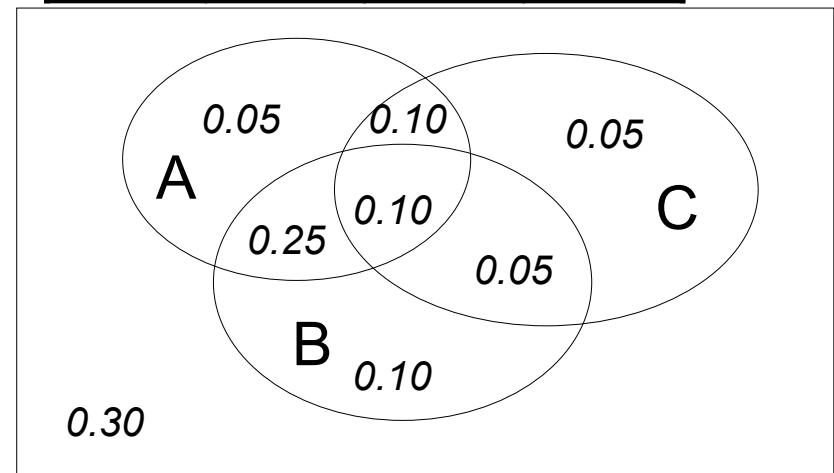instead of  h: X $\rightarrow$   Y,
learn          P(Y | X)

# Discriminative vs Generative Models

| Discriminative Models | Generative Models |
|---|---|
| Directly learn the function mapping $$h: X \rightarrow y$$ **or,** Calculate posterior distribution $$P(y\|X)$$ | Calculate posterior distribution $$P(y\|X)$$ from $P(X\|y)$ and $P(y)$ <br><br> **But Joint Distribution** $$P(X, y) = P(X\|y)\, P(y)$$ |
| 1. Assume some functional form for $P(y\|X)$ <br> 2. Estimate parameters of $P(y\|X)$ directly from training data | |

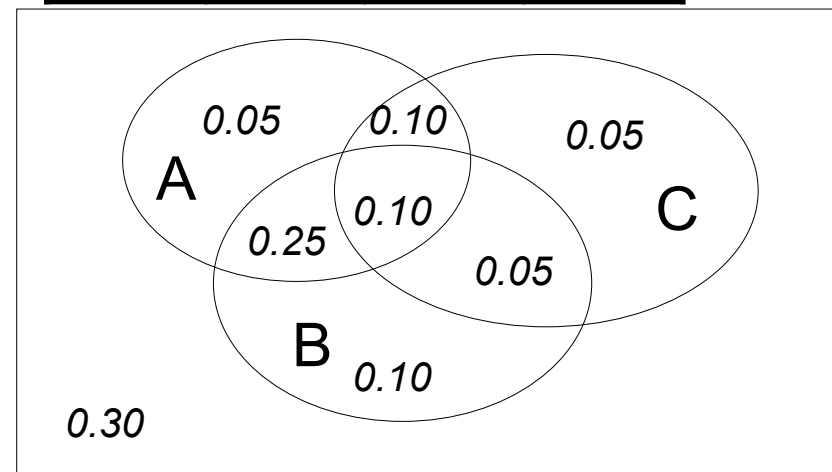# The Joint Distribution

*Example: Boolean variables A, B, C*

Recipe for making a joint distribution of M variables:

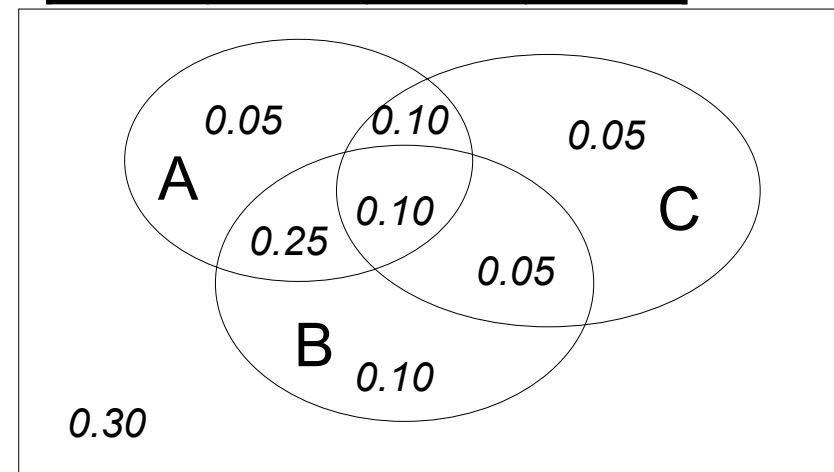| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | 0.10 |



[A. Moore]

# The Joint Distribution

*Example: Boolean variables A, B, C*

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values (M Boolean variables → $2^M$ rows).

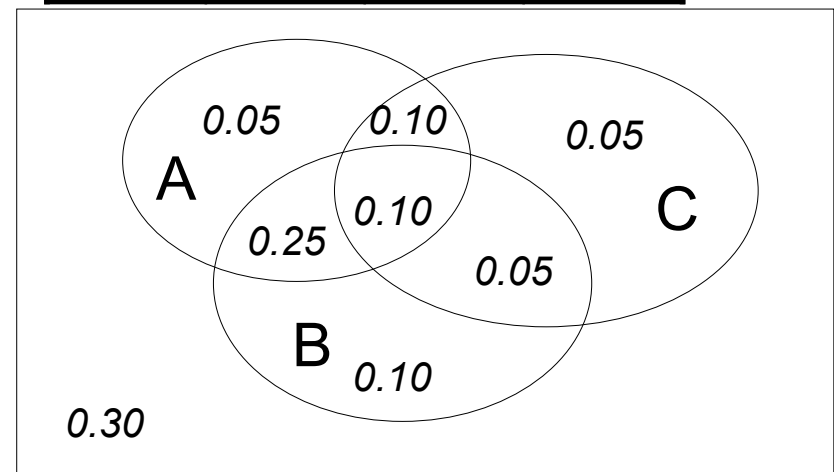| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | 0.10 |



[A. Moore]

# The Joint Distribution

*Example: Boolean variables A, B, C*

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values (M Boolean variables → $2^M$ rows).

2. For each combination of values, say how probable it is.

| A | B | 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 |



[A. Moore]

# The Joint Distribution

*Example: Boolean variables A, B, C*

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values (M Boolean variables → $2^M$ rows).

2. For each combination of values, say how probable it is.

3. If you subscribe to the axioms of probability, those probabilities must sum to 1.

| A | B | 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 |

[A. Moore]

# Using the Joint Distribution

| gender | hours_worked | wealth | | |
|--------|--------------|--------|--------|--------|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

One you have the JD you can ask for the probability of **any** logical expression involving these variables

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Using the Joint

| gender | hours_worked | wealth | | |
|--------|--------------|--------|-----------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

P(Poor Male) = 0.4654

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Using the Joint

| gender | hours_worked | wealth | | |
|--------|--------------|--------|---------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

P(Poor) = 0.7604

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Inference with the Joint



| gender | hours_worked | wealth | |
|--------|-------------|--------|----------|
| Female | v0:40.5- | poor | 0.253122 |
| | | rich | 0.0245895 |
| | v1:40.5+ | poor | 0.0421768 |
| | | rich | 0.0116293 |
| Male | v0:40.5- | poor | 0.331313 |
| | | rich | 0.0971295 |
| | v1:40.5+ | poor | 0.134106 |
| | | rich | 0.105933 |

$$P(E_1 \mid E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\displaystyle\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\displaystyle\sum_{\text{rows matching } E_2} P(\text{row})}$$

P(Male | Poor) = 0.4654 / 0.7604 = 0.612

# Learning and the Joint Distribution

| gender | hours_worked | wealth | | |
|--------|--------------|--------|-----------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

Suppose we want to learn the function f: <G, H> → W

Equivalently, P(W | G, H)

Solution: learn joint distribution from data, calculate P(W | G, H)

e.g., P(W=rich | G = female, H = 40.5- ) =

$$\frac{P(W=r \wedge G=f \wedge H=40\sim)}{P(\quad G=f \wedge H=40\sim)} = \frac{.024}{.277} \approx .09$$

[A. Moore]

sounds like the solution to learning $h: X \rightarrow Y$, or $P(Y \mid X)$.


Are we done?

sounds like the solution to
learning h: X → Y,
or P(Y | X).   $2^{10} = 1024$

Main problem: learning P(Y|X)
can require more data than we have

consider learning Joint Dist. with 100 attributes
# of rows in this table? $2^{100} \gtrapprox 1000^{16} = 10^{30}$
# of people on earth?
fraction of rows with 0 training examples? 0.9999

# What to do?

1.  Be smart about how we estimate probabilities from sparse data
    -   maximum likelihood estimates
    -   maximum a posteriori estimates


2.  Be smart about how to represent joint distributions
    -   Bayes networks, graphical models

# 1. Be smart about how we estimate probabilities

# Estimating Probability of Heads

X=1     X=0

- I show you the above coin $X$, and hire you to estimate the probability that it will turn up heads $(X = 1)$ or tails $(X = 0)$

- You flip it repeatedly, observing

    - it turns up heads $\alpha_1$ times
    - it turns up tails $\alpha_0$ times

- Your estimate for $P(X = 1)$ is....?

$$P(x=1) \approx \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

# Estimating $\theta = P(X=1)$

X=1   X=0

Test A:

$\alpha_1$              $\alpha_0$

100 flips: 51 Heads (X=1), 49 Tails (X=0)

$$\frac{\alpha_1}{\alpha_1 + \alpha_0} = \frac{51}{100} \rightarrow \hat{P}(X=1) = 0.51$$

Test B:

$\alpha_1$              $\alpha_0$

3 flips:  2 Heads (X=1), 1 Tails (X=0)

$$= \frac{2}{2+1} = 0.666$$

# Estimating θ = P(X=1)

X=1    X=0

Case C: (online learning)

- keep flipping, want single learning algorithm that gives reasonable estimate after each flip

$$\alpha_1 = \text{\# obs. heads } (x=1)$$

$$\alpha_0 = \text{\# obs } X=0$$

$$\beta_1 = \text{\# hallucinated } X=1\text{'s}$$

$$\beta_0 = \text{\# hallucinated } X=0\text{'s}$$

$$n = \alpha_1 + \alpha_0$$

$$\frac{\alpha_1 + 10}{(\alpha_1 + 10) + (\alpha_0 + 10)} \rightarrow \frac{(\alpha_1 + \beta_1)}{(\alpha_1 + \beta_1) + (\alpha_0 + \beta_0)}$$

# Principles for Estimating Probabilities

Principle 1 (maximum likelihood):

- choose parameters $\theta$ that maximize **P(data | $\theta$)**

- e.g., $$\hat{\theta}^{MLE} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

$$\frac{P(data|\theta)\,P(\theta)}{P(data)}$$

$$=$$

Principle 2 (maximum a posteriori prob.):

- choose parameters $\theta$ that maximize **P($\theta$ | data)**

- e.g.

$$\hat{\theta}^{MAP} = \frac{\alpha_1 + \#\text{hallucinated\_1s}}{(\alpha_1 + \#\text{hallucinated\_1s}) + (\alpha_0 + \#\text{hallucinated\_0s})}$$

# Maximum Likelihood Estimation

$P(X=1) = \theta$     $P(X=0) = (1-\theta)$

X=1    X=0

Data D: $= \{1 \quad 0 \quad 0 \quad 1\} \quad 1$

$P(D|\theta) = \theta \cdot (1-\theta) \cdot (1-\theta) \cdot \theta \cdot \theta = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$

Flips produce data D with $\alpha_1$ heads, $\alpha_0$ tails

- flips are independent, identically distributed 1's and 0's (Bernoulli)
- $\alpha_1$ and $\alpha_0$ are counts that sum these outcomes (Binomial)

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$$

# Maximum Likelihood Estimate for $\Theta$

$$\hat{\theta} = \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta)$$

$$= \arg\max_{\theta} \ln \theta^{\alpha_H}(1-\theta)^{\alpha_T}$$

- Set derivative to zero: $\boxed{\dfrac{d}{d\theta} \ln P(\mathcal{D} \mid \theta) = 0}$

[C. Guestrin]

$$\hat{\theta} = \arg\max_{\theta} \ln P(D|\theta)$$

■ Set derivative to zero: $\boxed{\dfrac{d}{d\theta} \ln P(\mathcal{D} \mid \theta) = 0}$

$$= \arg\max_{\theta} \ln \left[ \theta^{\alpha_1} (1 - \theta)^{\alpha_0} \right]$$

hint: $\dfrac{\partial \ln \theta}{\partial \theta} = \dfrac{1}{\theta}$

$$\frac{\partial}{\partial \theta} \quad \alpha_1 \ln \theta + \alpha_0 \ln(1-\theta)$$

$$\alpha_1 \frac{1}{\theta} + \alpha_0 \frac{\partial \ln(1-\theta)}{\partial \theta}$$

$$\boxed{0 = \alpha_1 \frac{1}{\theta} - \frac{\alpha_0}{1-\theta}}$$

$$\frac{\partial \ln(1-\theta)}{\partial(1-\theta)} \cdot \frac{\partial(1-\theta)}{\partial \theta}$$

$$\frac{1}{1-\theta} \qquad -1$$

$$\theta = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

# Summary:
## Maximum Likelihood Estimate

X=1    X=0

$P(X=1) = \theta$
$P(X=0) = 1\text{-}\theta$
(Bernoulli)

- Each flip yields boolean value for $X$

$$X \sim \text{Bernoulli: } P(X) = \theta^X (1 - \theta)^{(1-X)}$$

- Data set $D$ of independent, identically distributed (iid) flips produces $\alpha_1$ ones, $\alpha_0$ zeros (Binomial)

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1}(1 - \theta)^{\alpha_0}$$

$$\hat{\theta}^{MLE} = \text{argmax}_\theta \, P(D|\theta) = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

# Principles for Estimating Probabilities

Principle 1 (maximum likelihood):

- choose parameters $\theta$ that maximize $P(\text{data} \mid \theta)$

Principle 2 (maximum a posteriori prob.):

- choose parameters $\theta$ that maximize

$$P(\theta \mid \text{data}) = \frac{P(\text{data} \mid \theta)\, P(\theta)}{P(\text{data})}$$

# Beta prior distribution – P(θ)

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

- Likelihood function: $P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1 - \theta)^{\alpha_T}$
- Posterior: $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

# Beta prior distribution – P($\theta$)

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

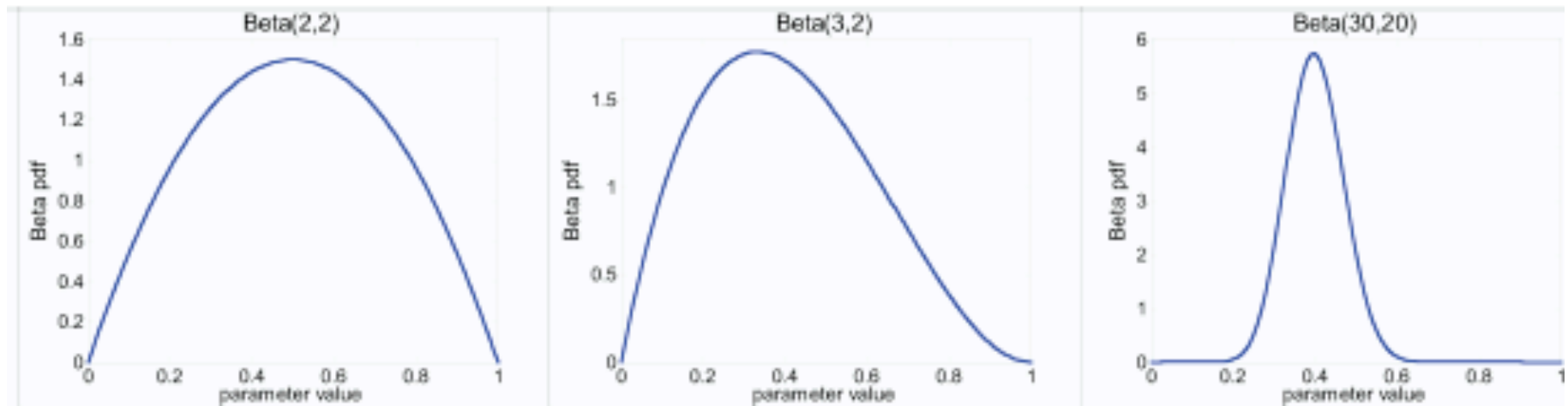- **Likelihood function:** $P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1 - \theta)^{\alpha_T}$

- **Posterior:** $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

$$\propto \theta^{\alpha_H + \beta_H - 1}(1-\theta)^{\alpha_T + \beta_T - 1}$$

$$\hat{\theta}^{MAP} = \frac{(\alpha_H + \beta_H - 1)}{(\alpha_H + \beta_H - 1) + (\alpha_T + \beta_T - 1)}$$

# Beta prior distribution – P(θ)

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$



[C. Guestrin]

## Eg. 1 Coin flip problem

Likelihood is ~ Binomial

$$P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1 - \theta)^{\alpha_T}$$

If prior is Beta distribution,

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

Then posterior is Beta distribution

$$P(\theta|D) \sim Beta(\alpha_H + \beta_H, \alpha_H + \beta_H)$$

and MAP estimate is therefore

$$\hat{\theta}^{MAP} = \frac{\alpha_H + \beta_H - 1}{(\alpha_H + \beta_H - 1) + (\alpha_T + \beta_T - 1)}$$

Eg. 2 Dice roll problem (6 outcomes instead of 2)

Likelihood is ~ Multinomial($\theta = \{\theta_1, \theta_2, \dots, \theta_k\}$)

$$P(\mathcal{D} \mid \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \dots \theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\theta_1^{\beta_1 - 1} \theta_2^{\beta_2 - 1} \dots \theta_k^{\beta_k - 1}}{B(\beta_1, \dots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \dots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta|D) \sim \text{Dirichlet}(\beta_1 + \alpha_1, \dots, \beta_k + \alpha_k)$$

and MAP estimate is therefore

$$\hat{\theta}_i^{MAP} = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

# MLE vs MAP

- Goal is same
  - Calculate posterior distribution
  - Probability of Data, $P(y|X)$
- **MLE**
  - **Does not use prior**
  - Starts with an assumption
- **MAP**
  - Uses Bayes Rule
  - Uses Prior
  - $P(y|X) \propto P(X|y)P(y)$

# Some terminology

- Likelihood function:  P(data | θ)
- Prior: P(θ)
- Posterior: P(θ | data)

- Conjugate prior: P(θ) is the conjugate prior for likelihood function P(data | θ) if the forms of P(θ) and P(θ | data) are the same.

# You should know

- Probability basics
  - random variables, conditional probs, …
  - Bayes rule
  - Joint probability distributions
  - calculating probabilities from the joint distribution
- Estimating parameters from data
  - maximum likelihood estimates
  - maximum a posteriori estimates
  - distributions – binomial, Beta, Dirichlet, …
  - conjugate priors

# Extra slides

# Independent Events

- Definition: two events A and B are *independent* if $P(A \wedge B) = P(A) * P(B)$

- Intuition: knowing A tells us nothing about the value of B (and vice versa)

# Picture "A independent of B"

# Expected values

Given a discrete random variable X, the expected value
of X, written E[X] is

$$E[X] = \sum_{x \in \mathcal{X}} x P(X = x)$$

Example:

| X | P(X) |
|---|------|
| 0 | 0.3 |
| 1 | 0.2 |
| 2 | 0.5 |

# Expected values

Given discrete random variable X, the expected value of
X, written E[X] is

$$E[X] = \sum_{x \in \mathcal{X}} x P(X = x)$$

We also can talk about the expected value of functions
of X

$$E[f(X)] = \sum_{x \in \mathcal{X}} f(x) P(X = x)$$

# Covariance

Given two discrete r.v.'s X and Y, we define the covariance of X and Y as

$$Cov(X, Y) = E[(X - E(X))(Y - E(Y))]$$

e.g., X=gender, Y=playsFootball

or    X=gender, Y=leftHanded

Remember: $E[X] = \sum_{x \in \mathcal{X}} x P(X = x)$