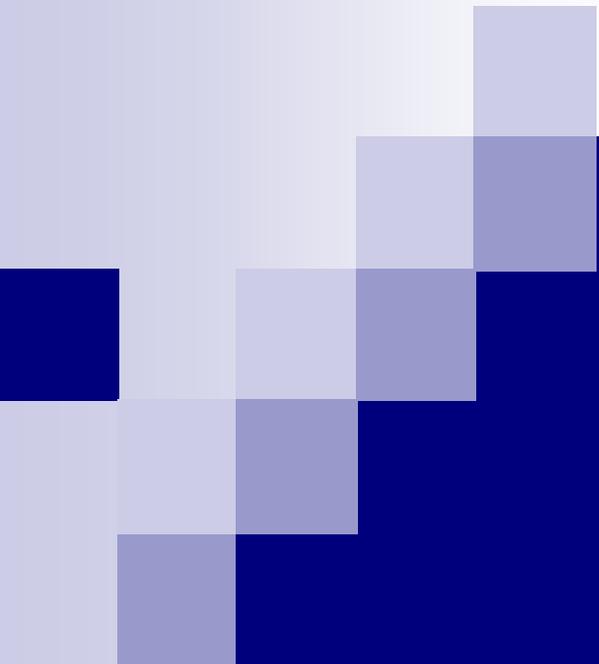# What's learning?
# Point Estimation

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

January 18th, 2005

# Growth of Machine Learning

- Machine learning is preferred approach to
  - Speech recognition, Natural language processing
  - Computer vision
  - Medical outcomes analysis
  - Robot control
  - …
- This trend is accelerating
  - Improved machine learning algorithms
  - Improved data capture, networking, faster computers
  - Software too complex to write by hand
  - New sensors / IO devices
  - Demand for self-customization to user, environment

# Syllabus

- Covers a wide range of Machine Learning techniques – from basic to state-of-the-art

- You will learn about the methods you heard about:
  - Naïve Bayes, logistic regression, nearest-neighbor, decision trees, boosting, neural nets, overfitting, regularization, dimensionality reduction, PCA, error bounds, VC dimension, SVMs, kernels, margin bounds, K-means, EM, mixture models, semi-supervised learning, HMMs, graphical models, active learning, reinforcement learning…

- Covers algorithms, theory and applications

- **It's going to be fun and hard work** ☺

# Prerequisites

- Probabilities
  - Distributions, densities, marginalization…
- Basic statistics
  - Moments, typical distributions, regression…
- Algorithms
  - Dynamic programming, basic data structures, complexity…
- Programming
  - Mostly your choice of language, but Matlab will be very useful
- We provide some background, but the class will be fast paced

- Ability to deal with "abstract mathematical concepts"

# Review Sessions

- Very useful!
  - Review material
  - Present background
  - Answer questions
- Thursdays, 5:00-6:30 in Wean Hall 5409
- First recitation is **tomorrow**
  - Review of probabilities
- Special recitation on Matlab
  - Jan. 25 Wed. 5:00-7:00pm, NSH 3305

# Staff

- Four Great TAs: Great resource for learning, interact with them!
  - Anton Chechetka, antonc@cs
  - Stanislav Funiak, sfuniak@cs
  - Andreas Krause, krausea@cs
  - Jure Leskovec, jure@cs

- Course General Czar
  - Terrill L. Frantz, TerrillFrantz@cmu
- Administrative Assistant
  - Monica Hopes, x8-5527, meh@cs

# First Point of Contact for HWs

- To facilitate interaction, a TA will be assigned to each homework question – This will be your "first point of contact" for this question
  - But, you can always ask any of us

- For e-mailing instructors, always use:
  - 10701-instructors@cs.cmu.edu

- For announcements, subscribe to:
  - 10701-announce@cs
  - https://mailman.srv.cs.cmu.edu/mailman/listinfo/10701-announce

# All Text Books are Optional, *but very useful*

- *Machine Learning*, Tom Mitchell
- *Pattern Classification (2nd Edition)*, Duda, Hart and Stork
- *Neural Networks for Pattern Recognition*, Chris Bishop

# Grading

- 5 homeworks (30%)
  - First one goes out 1/23
- Final project (20%)
  - Details out March 1$^{st}$
- Midterm (20%)
  - March 8$^{th}$
- Final (30%)
  - TBD by registrar

# Homeworks

- Homeworks are hard, start early ☺
- Due in the beginning of class
- 3 late days for the semester
- After late days are used up:
    - Half credit within 48 hours
    - Zero credit after 48 hours
- All homeworks **must be handed in**, even for zero credit
- Late homeworks handed in to Monica Hopes, WEH 4616

- Collaboration
    - You may **discuss** the questions
    - Each student writes their own answers
    - Write on your homework anyone with whom you collaborate

# Enjoy!

- ML is becoming ubiquitous in science, engineering and beyond
- This class should give you the basic foundation for applying ML and developing new methods
- The fun begins…

# What is Machine Learning ?

# Machine Learning

Study of algorithms that
- improve their <u>performance</u>
- at some <u>task</u>
- with <u>experience</u>

# Object detection

(Prof. H. Schneiderman)

Example training images
for each orientation

# Text classification

webpage

$X \longrightarrow \{C, P, U, ... \}$



→ Company home page

 vs

→ Personal home page

 vs

Univeristy home page

 vs

…

Reading
a noun
(vs verb)

[Rustandi et al., 2005]
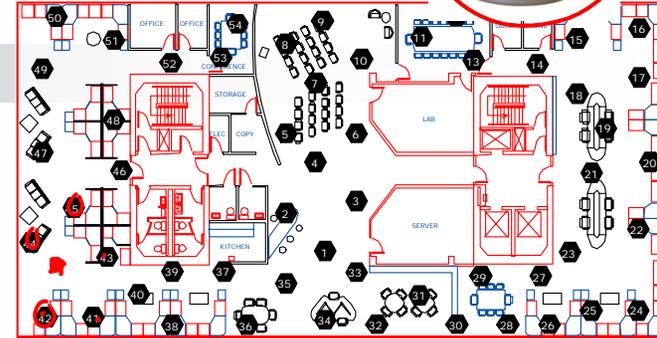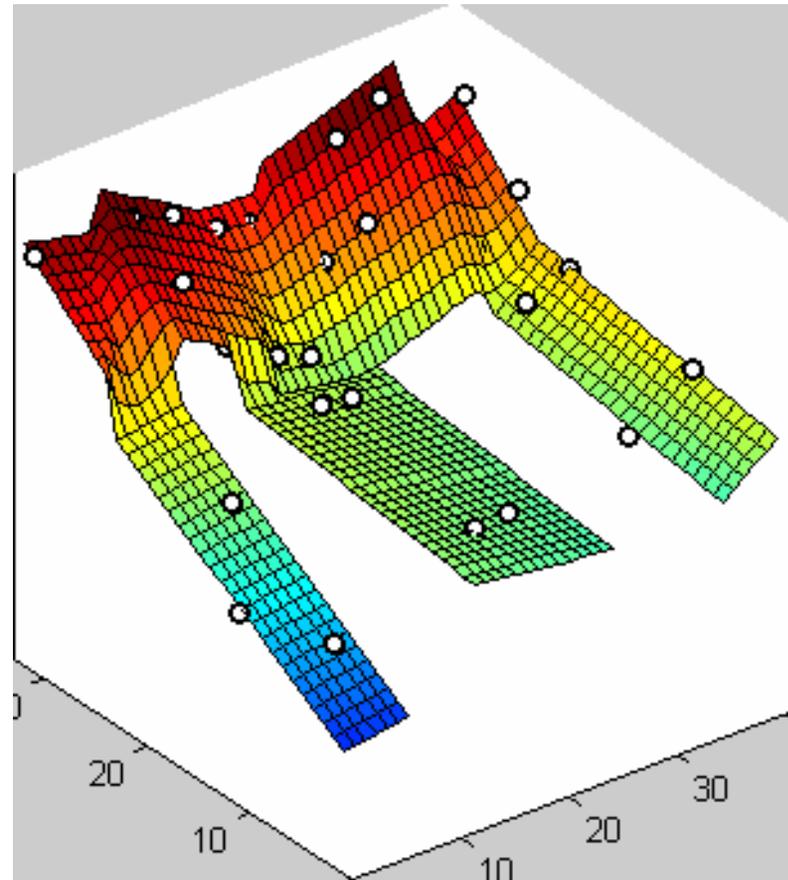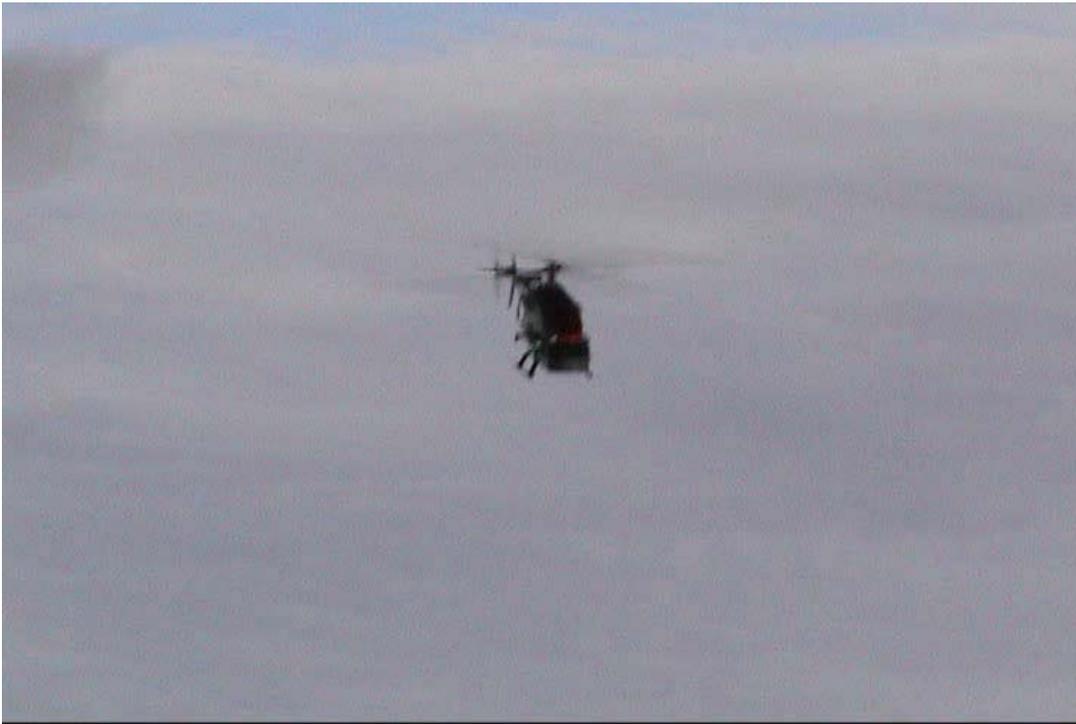
# Modeling sensor data



- Measure temperatures at some locations
- Predict temperatures throughout the environment

[Guestrin et al. '04]

# Learning to act



[Ng et al. '05]

- **Reinforcement learning**
- **An agent**
  - ☐ Makes sensor observations
  - ☐ Must select action
  - ☐ Receives rewards
    - positive for "good" states
    - negative for "bad" states

# Your first consulting job

- A billionaire from the suburbs of Seattle asks you a question:
  - He says: I have thumbtack, if I flip it, what's the probability it will fall with the nail up?
  - You say: Please flip it a few times:

  - You say: The probability is: $\frac{3}{5}$
  - **He says: Why???**
  - You say: Because…

# Thumbtack – Binomial Distribution

■ P(Heads) = $\theta$,  P(Tails) = 1-$\theta$

$\alpha_H = 3, \alpha_T = 2$

$P(HTTHH) = \theta(1-\theta)(1-\theta)\theta\theta = \theta^3(1-\theta)^2$

■ Flips are i.i.d.:

  ☐ Independent events

  ☐ Identically distributed according to Binomial distribution

■ Sequence $D$ of $\alpha_H$ Heads and $\alpha_T$ Tails

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

# Maximum Likelihood Estimation

- **<u>Data</u>:** Observed set $D$ of $\alpha_H$ Heads and $\alpha_T$ Tails
- **Hypothesis:** Binomial distribution
- Learning $\theta$ is an optimization problem
  - ☐ What's the objective function?

$D = (HTTHH) \rightarrow$ want pick a coin $\theta$ that would generate $D$

- <u>MLE</u>: Choose $\theta$ that maximizes the probability of observed data:   Pick coin $\theta$ max. $P(HTTHH|\theta)$

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

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

# Your first learning algorithm

$\ln ab = \ln a + \ln b; \quad \ln a^b = b \ln a; \frac{d}{d\theta}\ln\theta = \frac{1}{\theta}$

$\frac{d}{d\theta}\ln(1-\theta) = \frac{-1}{1-\theta}$

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

estimate

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

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

$\dfrac{d}{d\theta}\ln\left[\theta^{\alpha_H}(1-\theta)^{\alpha_T}\right]$

$= \dfrac{d}{d\theta}\left[\ln\theta^{\alpha_H} + \ln(1-\theta)^{\alpha_T}\right]$

$= \dfrac{d}{d\theta}\left[\alpha_H \ln\theta + \alpha_T \ln(1-\theta)\right]$

$= \dfrac{\alpha_H}{\theta} - \dfrac{\alpha_T}{1-\theta} = 0$

$\theta = \dfrac{\alpha_H}{\alpha_H + \alpha_T}$

$\alpha_H = 3$

$\alpha_T = 2$

$\theta = \dfrac{3}{5}\ \text{!!} \smile$

# How many flips do I need?

$$\widehat{\theta} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

- Billionaire says: I flipped 3 heads and 2 tails.  $\theta = \frac{3}{5}$
- You say: $\theta$ = 3/5, I can prove it!
- He says: What if I flipped 30 heads and 20 tails?
- You say: Same answer, I can prove it!
- **He says: What's better?**
- You say: Humm… The more the merrier???
- He says: Is this why I am paying you the big bucks???

# Simple bound (based on Hoeffding's inequality)

- For $N = \alpha_H + \alpha_T$, and $\widehat{\theta} = \dfrac{\alpha_H}{\alpha_H + \alpha_T}$

  5
  50
  $\vdots$

- Let $\theta^*$ be the true parameter, for any $\varepsilon > 0$:

$$P(\,|\,\widehat{\theta} - \theta^*\,| \geq \epsilon\,) \;\leq\; 2e^{-2N\epsilon^2}$$

0.1

# PAC Learning

- PAC: (Probably) Approximate Correct

- Billionaire says: I want to know the thumbtack parameter $\theta$, within $\varepsilon = 0.1$, with probability at least $1-\delta = 0.95$. How many flips?

$$P(|\,\widehat{\theta} - \theta^*\,| \geq \underset{0.1}{\epsilon}) \;\leq\; 2e^{-2N\epsilon^2} \leq \delta = 0.05$$

$$2e^{-2N\varepsilon^2} \leq \delta$$

$$\ln\{2e^{-2N\varepsilon^2}\} \leq \ln\delta$$

$$\ln 2 - 2N\varepsilon^2 \leq \ln\delta$$

$$N \geq \frac{1}{2\varepsilon^2}\left[\ln 2 - \ln\delta\right]$$

$$N \geq \frac{1}{2\varepsilon^2}\ln\frac{2}{\delta}$$

$$N \geq \frac{1}{2(0.1)^2}\ln\frac{2}{0.05}$$

# What about prior

- Billionaire says: Wait, I know that the thumbtack is "close" to 50-50. What can you? *say*

- **You say: I can learn it the Bayesian way…**

- Rather than estimating a single $\theta$, we obtain a distribution over possible values of $\theta$

# Bayesian Learning

MLE: $\text{argmax}_{\theta} P(D|\theta)$

- Use Bayes rule:

$$\underbrace{P(\theta \mid \mathcal{D})}_{\text{Posterior}} = \frac{\overbrace{P(\mathcal{D} \mid \theta)}^{\text{likelihood}} \overbrace{P(\theta)}^{\text{prior}}}{\underbrace{P(\mathcal{D})}_{\substack{\text{Normalization} \\ \text{constant}}}}$$

- Or equivalently:

$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta) P(\theta)$$

# Bayesian Learning for Thumbtack

$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$$

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

- What about prior?
  - Represent expert knowledge
  - Simple posterior form
- Conjugate priors:
  - Closed-form representation of posterior
  - **For Binomial, conjugate prior is Beta distribution**
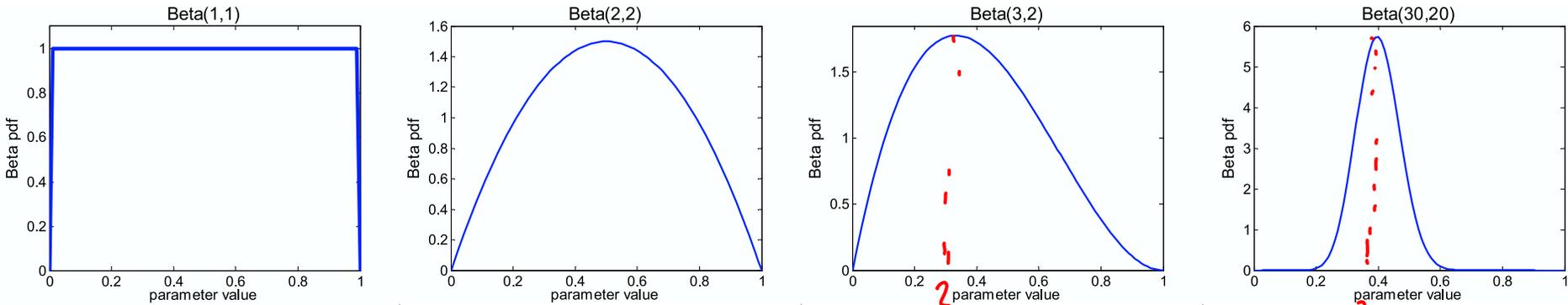
$P(\theta \mid D) \propto P(D \mid \theta) P(\theta)$

likelihood $\rightarrow$ multinomial
prior $\rightarrow$ Beta
postior $\rightarrow$ Closed form
(Beta)

# Beta prior distribution – P(θ)

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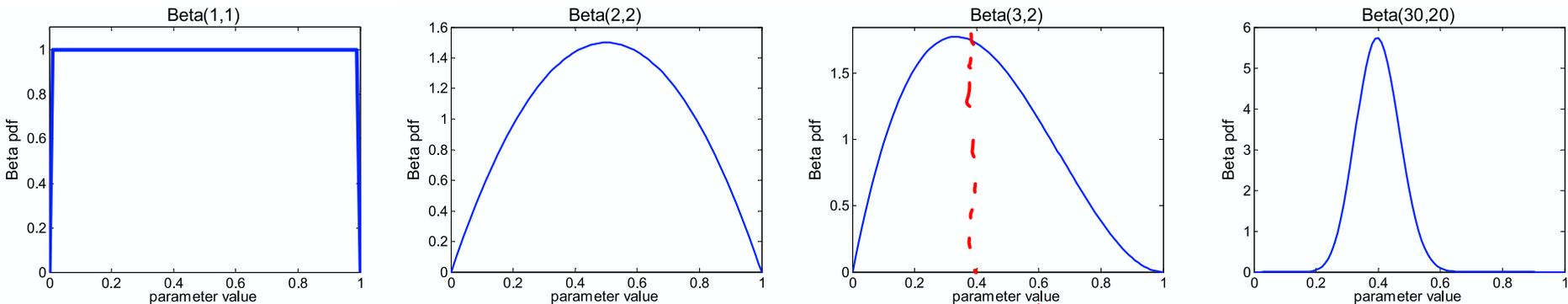


Beta(1,1)  Beta(2,2)  Beta(3,2)  Beta(30,20)

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

# Posterior distribution

- Prior: $Beta(\beta_H, \beta_T)$

- Data: $\alpha_H$ heads and $\alpha_T$ tails

- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

start with
$\beta_H = 1, \beta_T = 1$

$D : \alpha_H = 2$
$\alpha_T = 1$



Beta(1,1)    Beta(2,2)    Beta(3,2)    Beta(30,20)

start

$(1-\theta) = 2/5$
end up

# Using Bayesian posterior


Beta(30,20)

- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- Bayesian inference:
  - ☐ No longer single parameter:

$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta \mid \mathcal{D}) d\theta$$

earnings   posterior

  - ☐ Integral is often hard to compute

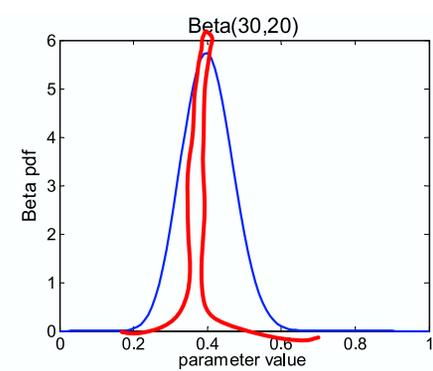# MAP: Maximum a posteriori approximation

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$



Beta(30,20)

*mode is MAP*

$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta \mid \mathcal{D}) d\theta$$
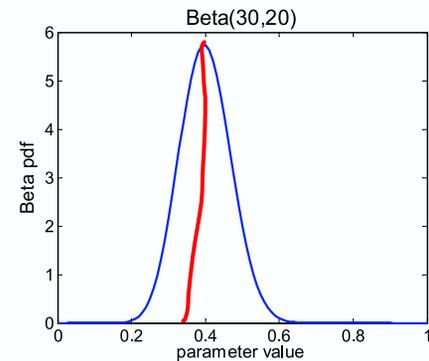
- As more data is observed, Beta is more certain

- MAP: use most likely parameter:

$$\widehat{\theta} = \arg\max_{\theta} P(\theta \mid \mathcal{D}) \qquad E[f(\theta)] \approx f(\widehat{\theta})$$

# MAP for Beta distribution


Beta(30,20)

$$P(\theta \mid \mathcal{D}) = \frac{\theta^{\beta_H + \alpha_H - 1}(1 - \theta)^{\beta_T + \alpha_T - 1}}{B(\beta_H + \alpha_H, \beta_T + \alpha_T)} \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

*multinomial ~ like*

$\alpha_H = 3$
$\alpha_T = 2$
$\beta_H, \beta_T$ extra data

- ■ MAP: use most likely parameter:

$$\widehat{\theta} = \arg\max_{\theta} P(\theta \mid \mathcal{D}) = \frac{\beta_H + \alpha_H - 1}{\beta_H + \alpha_H + \beta_T + \alpha_T - 2}$$

- ■ Beta prior equivalent to extra thumbtack flips
- ■ As $N \to \infty$, prior is "forgotten"
- ■ **But, for small sample size, prior is important!**

# What you need to know

- Go to the recitation on intro to probabilities
  - And, other recitations too
- Point estimation:
  - MLE
  - Bayesian learning
  - MAP