

Linear Regression

Machine Learning – CSE546
 Sham Kakade
 University of Washington
 Oct 4, 2016

©2016 Sham Kakade

Announcements: *Write your own code.*

- TA office hours posted on website
- Recitation this week: Python
- HW1 posted
- Addcodes: decided by tomorrow. You will be contacted by email.

- Today:
 - MLE continued
 - Regression

©2016 Sham Kakade

What about continuous variables?

- Billionaire says: If I am measuring a continuous variable, what can you do for me?
- **You say: Let me tell you about Gaussians...**

$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

©2016 Sham Kakade

Some properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
 - $X \sim N(\mu, \sigma^2)$
 - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$
- Sum of Gaussians
 - $X \sim N(\mu_X, \sigma_X^2)$
 - $Y \sim N(\mu_Y, \sigma_Y^2)$
 - $Z = X + Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$

©2016 Sham Kakade

Learning a Gaussian

- Collect a bunch of data
 - Hopefully, i.i.d. samples
 - e.g., exam scores

- Learn parameters

- Mean
- Variance

$$P(x | \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

©2016 Sham Kakade

5

$x_1 = 98$
 \vdots
 $x_N = 32$ also,
 $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i$ ← MLE.
 (why?)

MLE for Gaussian

- Prob. of i.i.d. samples $D = \{x_1, \dots, x_N\}$:

$$P(D | \mu, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

$$\hat{\mu}_{MLE}, \hat{\sigma}_{MLE} = \arg \max_{\mu, \sigma} \ln P(D | \mu, \sigma)$$

- Log-likelihood of data:

$$\begin{aligned} \ln P(D | \mu, \sigma) &= \ln \left[\left(\frac{1}{\sigma\sqrt{2\pi}} \right)^N \prod_{i=1}^N e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right] \\ &= -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

©2016 Sham Kakade

6

$P_D(x_1, \dots, x_N | \mu, \sigma)$
 $= P(x_1 | \mu, \sigma) \dots P(x_N | \mu, \sigma)$
 by i.i.d.

Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

$$\frac{d}{d\mu} \ln P(D | \mu, \sigma) = \frac{d}{d\mu} \left[-N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$- \sum_{i=1}^N \frac{d}{d\mu} \frac{(x_i - \mu)^2}{2\sigma^2} = - \sum_{i=1}^N \frac{-2(x_i - \mu)}{2\sigma^2} = 0$$

$$\hat{\mu} = \frac{\sum_{i=1}^N x_i}{N}$$

©2016 Sham Kakade

7

MLE for variance

- Again, set derivative to zero:

$$\begin{aligned} \frac{d}{d\sigma} \ln P(D | \mu, \sigma) &= \frac{d}{d\sigma} \left[-N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= \frac{d}{d\sigma} \left[-N \ln \sigma\sqrt{2\pi} \right] - \sum_{i=1}^N \frac{d}{d\sigma} \left[\frac{(x_i - \mu)^2}{2\sigma^2} \right] \end{aligned}$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$

©2016 Sham Kakade

8

Learning Gaussian parameters

- MLE:

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

- BTW. MLE for the variance of a Gaussian is **biased**
 - Expected result of estimation is **not** true parameter!
 - Unbiased variance estimator:

$$\hat{\sigma}_{unbiased}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2$$

©2016 Sham Kabade

Prediction of continuous variables

- Billionaire says: Wait, that's not what I meant!
- You say: Chill out, dude.
- She says: I want to predict a continuous variable for continuous inputs: I want to predict salaries from GPA.
- You say: **I can regress that...**

©2016 Sham Kabade

10

The regression problem

- Instances: $\langle \mathbf{x}_j, t_j \rangle$
- Learn: Mapping from \mathbf{x} to $t(\mathbf{x})$

Hypothesis space:

- Given, basis functions $H = \{h_1, \dots, h_K\}$
- Find coeffs $\mathbf{w} = \{w_1, \dots, w_K\}$ $t(\mathbf{x}) \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$
- Why is this called linear regression???

- model is linear in the parameters

- Precisely, minimize the **residual squared error**:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{j=1}^N \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

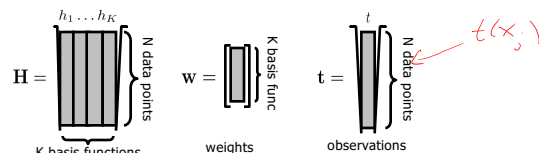
©2016 Sham Kabade

11

The regression problem in matrix notation

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \underbrace{(\mathbf{H}\mathbf{w} - \mathbf{t})^T (\mathbf{H}\mathbf{w} - \mathbf{t})}_{\text{residual error}}$$



©2016 Sham Kabade

12

Minimizing the Residual

scalar case $\frac{d}{dw} (w^T H w - t^T H w)$

$w^* = \arg \min_w \underbrace{(Hw - t)^T (Hw - t)}_{\text{residual error}} = 2\alpha (w - t)$

$\nabla_w F(w) = 0$

$= 2H^T (Hw - t) = 0$

$H^T H w - H^T t = 0$

$w^* = (H^T H)^{-1} H^T t$

in matrix case $\nabla_w (Hw - t)^T (Hw - t)$

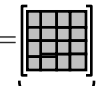
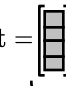
$= 2H^T (Hw - t)$

©2016 Sham Kakade 13

Regression solution = simple matrix operations

$w^* = \arg \min_w \underbrace{(Hw - t)^T (Hw - t)}_{\text{residual error}}$

solution: $w^* = \underbrace{(H^T H)^{-1}}_{A^{-1}} \underbrace{H^T t}_b = \mathbf{A}^{-1} \mathbf{b}$

where $\mathbf{A} = H^T H =$  $\mathbf{b} = H^T t =$ 

$k \times k$ matrix for k basis functions $k \times 1$ vector

©2016 Sham Kakade 14

But, why?

- Billionaire again, she says: Why sum squared error???
- You say: Gaussians, Gaussians... *Assume $\epsilon_x \sim \mathcal{N}(0, \sigma^2)$*
- Model: prediction is linear function plus Gaussian noise
 - $t(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}) + \epsilon_x$ • $P_\sigma(t(\mathbf{x}) | \mathbf{x}) = \mathcal{N}(\sum_i w_i h_i(\mathbf{x}), \sigma^2)$
 - $t(\mathbf{x})$ i.i.d.
- Learn w using MLE *one point*

$$P(t | \mathbf{x}, w, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{[t - \sum_i w_i h_i(\mathbf{x})]^2}{2\sigma^2}}$$

©2016 Sham Kakade 15

Maximizing log-likelihood

$P(t(\mathbf{x}_j), t(\mathbf{x}_j) | w, \sigma) \dots$

Maximize:

$$\ln P(\mathcal{D} | w, \sigma) = \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{j=1}^N e^{-\frac{[t_j - \sum_i w_i h_i(\mathbf{x}_j)]^2}{2\sigma^2}}$$

argmax w $\ln(\dots) - \frac{N}{2\sigma^2} (t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j))^2$

const

$= \arg \min_w \sum_{j=1}^N (t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j))^2$

Least-squares Linear Regression is MLE for Gaussians!!!

©2016 Sham Kakade 16

Bias-Variance Tradeoff

Machine Learning – CSE546
Sham Kakade
University of Washington
Oct 4, 2016

©2016 Sham Kakade 17

Bias-Variance tradeoff – Intuition

- Model too “simple” → does not fit the data well
 - A biased solution
- Model too complex → small changes to the data, solution changes a lot
 - A high-variance solution

©2016 Sham Kakade 18

(Squared) Bias of learner

- Given dataset D with N samples, learn function $h_D(x)$
- If you sample a different dataset D' with N samples, you will learn different $h_{D'}(x)$
- **Expected hypothesis:** $E_D[h_D(x)]$
- **Bias:** difference between what you expect to learn and truth
 - Measures how well you expect to represent true solution
 - Decreases with more complex model
 - Bias² at one point x : $(f(x) - \bar{h}_D(x))^2$
 - Average Bias²: $E_x \{ (f(x) - \bar{h}_D(x))^2 \}$

©2016 Sham Kakade 19

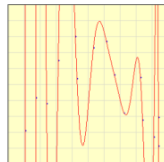
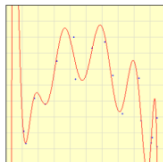
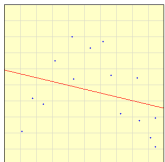
Variance of learner

- Given dataset D with N samples, learn function $h_D(x)$
- If you sample a different dataset D' with N samples, you will learn different $h_{D'}(x)$
- **Variance:** difference between what you expect to learn and what you learn from a particular dataset
 - Measures how sensitive learner is to specific dataset
 - Decreases with simpler model
 - Variance at one point x : $E_D [(h_D(x) - \bar{h}_D(x))^2]$
 - **Average variance:** $E_x E_D [(h_D(x) - \bar{h}_D(x))^2]$

©2016 Sham Kakade 20

Bias-Variance Tradeoff

- Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance



Select points by clicking on the graph or press [Example](#)

Degree of polynomial: FEY to X FEK to Y

[Calculate](#) [View Polynomial](#) [Reset](#)

©2016 Sham Kabade 21

Bias-Variance Decomposition of Error

$$\bar{h}_N(x) = E_D[h_D(x)]$$

- Expected mean squared error: $MSE = E_D \left[E_x \left[(t(x) - h_D(x))^2 \right] \right]$
- To simplify derivation, drop x :
- Expanding the square:

Moral of the Story: Bias-Variance Tradeoff Key in ML

- Error can be decomposed:

$$MSE = E_D \left[E_x \left[(t(x) - h_D(x))^2 \right] \right]$$

$$= E_x \left[(t(x) - \bar{h}_N(x))^2 \right] + E_D \left[E_x \left[(\bar{h}_N(x) - h_D(x))^2 \right] \right]$$
- Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance

What you need to know

- Regression
 - Basis function = features
 - Optimizing sum squared error
 - Relationship between regression and Gaussians
- Bias-variance trade-off
- Play with Applet

Overfitting

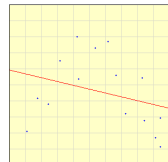
Machine Learning – CSE546
 Sham Kakade
 University of Washington
 Oct 4, 2016

©2016 Sham Kakade

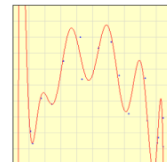
25

Bias-Variance Tradeoff

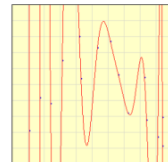
- Choice of hypothesis class introduces learning bias
 - More complex class → less bias
 - More complex class → more variance



Select points by clicking on the graph or press [Example](#)
 Degree of polynomial: Fix to X Fix to Y



Select points by clicking on the graph or press [Example](#)
 Degree of polynomial: Fix to X Fix to Y



Select points by clicking on the graph or press [Example](#)
 Degree of polynomial: Fix to X Fix to Y

©2016 Sham Kakade

26

Training set error $w^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$

- Given a dataset (Training data)
- Choose a loss function
 - e.g., squared error (L_2) for regression
- **Training set error:** For a particular set of parameters, loss function on training data:

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2016 Sham Kakade

27

Training set error as a function of model complexity

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$



Select points by clicking on the graph or press [Example](#)
 Degree of polynomial: Fix to X Fix to Y



Select points by clicking on the graph or press [Example](#)
 Degree of polynomial: Fix to X Fix to Y

©2016 Sham Kakade

28

Prediction error

- Training set error can be poor measure of “quality” of solution
- **Prediction error:** We really care about error over all possible input points, not just training data:

$$\begin{aligned} error_{true}(\mathbf{w}) &= E_{\mathbf{x}} \left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 \right] \\ &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

©2016 Sham Kakade

29

Prediction error as a function of model complexity

$$\begin{aligned} error_{train}(\mathbf{w}) &= \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2 \\ error_{true}(\mathbf{w}) &= \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x} \end{aligned}$$



©2016 Sham Kakade

30

Computing prediction error

- Computing prediction
 - Hard integral
 - May not know $t(\mathbf{x})$ for every \mathbf{x}

$$error_{true}(\mathbf{w}) = \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}$$

- Monte Carlo integration (sampling approximation)
 - Sample a set of i.i.d. points $\{\mathbf{x}_1, \dots, \mathbf{x}_M\}$ from $p(\mathbf{x})$
 - Approximate integral with sample average

$$error_{true}(\mathbf{w}) \approx \frac{1}{M} \sum_{j=1}^M \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2016 Sham Kakade

31

Why training set error doesn't approximate prediction error?

- Sampling approximation of prediction error:

$$error_{true}(\mathbf{w}) \approx \frac{1}{M} \sum_{j=1}^M \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

- Training error :

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

- Very similar equations!!!
 - Why is training set a bad measure of prediction error???

©2016 Sham Kakade

32

Why training set error doesn't approximate prediction error?

Because you cheated!!!

Training error good estimate for a single \mathbf{w} ,
But you optimized \mathbf{w} with respect to the training error,
and found \mathbf{w} that is good for this set of samples

Training error is a (optimistically) biased estimate of prediction error

- Very similar equations!!!
 - Why is training set a bad measure of prediction error???

©2016 Sham Kakade

33

Test set error

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_j \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

- Given a dataset, **randomly** split it into two parts:
 - Training data – $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{train}}}\}$
 - Test data – $\{\mathbf{x}_1, \dots, \mathbf{x}_{N_{\text{test}}}\}$
- Use training data to optimize parameters \mathbf{w}
- Test set error:** For the **final output** $\hat{\mathbf{w}}$, evaluate the error using:

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2016 Sham Kakade

34

Test set error as a function of model complexity

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$



©2016 Sham Kakade

35

Overfitting

- Overfitting:** a learning algorithm overfits the training data if it outputs a solution \mathbf{w} when there exists another solution \mathbf{w}' such that:

$$[error_{train}(\mathbf{w}) < error_{train}(\mathbf{w}')] \wedge [error_{true}(\mathbf{w}') < error_{true}(\mathbf{w})]$$

©2016 Sham Kakade

36

How many points to I use for training/testing?

- Very hard question to answer!
 - Too few training points, learned \mathbf{w} is bad
 - Too few test points, you never know if you reached a good solution

- Bounds, such as Hoeffding's inequality can help:

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

- More on this later this quarter, but still hard to answer
- Typically:
 - If you have a reasonable amount of data, pick test set "large enough" for a "reasonable" estimate of error, and use the rest for learning
 - If you have little data, then you need to pull out the big guns...
 - e.g., bootstrapping

©2016 Sham Kakade

37

Error estimators

$$error_{train}(\mathbf{w}) = \int_{\mathbf{x}} \left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x}) \right)^2 p(\mathbf{x}) d\mathbf{x}$$

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2016 Sham Kakade

38

Error as a function of number of training examples for a fixed model complexity

$$error_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

little data

infinite data

©2016 Sham Kakade

39

Error estimators

Be careful!!!

Test set only unbiased if you never ever ever ever do any any any any learning on the test data

For example, if you use the test set to select the degree of the polynomial... no longer unbiased!!!
(We will address this problem later in the quarter)

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left(t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

©2016 Sham Kakade

40

What you need to know

- True error, training error, test error
 - Never learn on the test data
 - Never learn on the test data
 - Never learn on the test data
 - Never learn on the test data
 - Never learn on the test data
- Overfitting

Bayesian Methods

Machine Learning – CSE546
Sham Kakade
University of Washington
Oct 4, 2016

©2016 Sham Kakade

45

What about prior

- Billionaire says: Wait, I know that the thumbtack is “close” to 50-50. What can you do for me now?
- **You say: I can learn it the Bayesian way...**
- Rather than estimating a single θ , we obtain a distribution over possible values of θ

©2016 Sham Kakade

46

Bayesian Learning

- Use Bayes rule:

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta)P(\theta)}{P(\mathcal{D})}$$

- Or equivalently:

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

©2016 Sham Kakade

47

Bayesian Learning for Thumbtack

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

- Likelihood function is simply Binomial:

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

- What about prior?
 - Represent expert knowledge
 - Simple posterior form
- Conjugate priors:
 - Closed-form representation of posterior
 - **For Binomial, conjugate prior is Beta distribution**

©2016 Sham Kakade

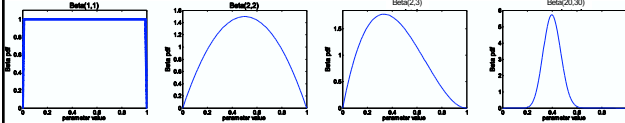
48

Beta prior distribution – $P(\theta)$

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

Mean:

Mode:



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

©2016 Sham Kabade

49

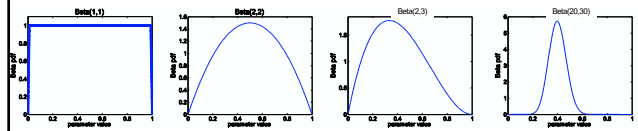
Posterior distribution

- Prior: $\text{Beta}(\beta_H, \beta_T)$

- Data: α_H heads and α_T tails

- Posterior distribution:

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



©2016 Sham Kabade

50

Using Bayesian posterior

- Posterior distribution:

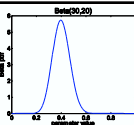
$$P(\theta | \mathcal{D}) \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- Bayesian inference:

- No longer single parameter:

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

- Integral is often hard to compute



©2016 Sham Kabade

51

MAP: Maximum a posteriori approximation

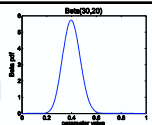
$$P(\theta | \mathcal{D}) \sim \text{Beta}(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

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

- As more data is observed, Beta is more certain

- MAP: use most likely parameter:

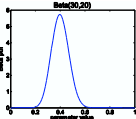
$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D}) \quad E[f(\theta)] \approx f(\hat{\theta})$$



©2016 Sham Kabade

52

MAP for Beta distribution



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

- MAP: use most likely parameter:

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

- Beta prior equivalent to extra thumbtack flips
- As $N \rightarrow 1$, prior is “forgotten”
- **But, for small sample size, prior is important!**