# Linear Regression

Machine Learning – CSE546

Sham Kakade

University of Washington

Oct 4, 2016

1

---

# 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…**

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

2

# Announcements:

- TA office hours posted on website
- Recitation this week: Python
- HW1 posted


- Today:
  - □ MLE continued
  - □ Regression

3

---

# Some properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
  - □ $X \sim N(\mu, \sigma^2)$
  - □ $Y = aX + b$ ➜ $Y \sim N(a\mu+b, a^2\sigma^2)$

- Sum of Gaussians
  - □ $X \sim N(\mu_X, \sigma^2_X)$
  - □ $Y \sim N(\mu_Y, \sigma^2_Y)$
  - □ $Z = X+Y$ ➜ $Z \sim N(\mu_X+\mu_Y, \sigma^2_X+\sigma^2_Y)$

4

# Learning a Gaussian

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

- Learn parameters
  - ☐ Mean
  - ☐ Variance

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

5

# MLE for Gaussian

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

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

- Log-likelihood of data:

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

$$= -N\ln\sigma\sqrt{2\pi} - \sum_{i=1}^{N} \frac{(x_i-\mu)^2}{2\sigma^2}$$

6

## Your second learning algorithm: MLE for mean of a Gaussian

■ What's MLE for mean?

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

7

## MLE for variance

■ Again, set derivative to zero:

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

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

8

# Learning Gaussian parameters

■ MLE:

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

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

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

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

# 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…**

# The regression problem

- **Instances:** $\langle \mathbf{x}_j, t_j \rangle$
- **Learn:** Mapping from x to t($\mathbf{x}$)
- **Hypothesis space:**
  - ☐ Given, basis functions

    $H = \{h_1, \ldots, h_K\}$
  - ☐ Find coeffs $\mathbf{w} = \{w_1, \ldots, w_k\}$

    $\underbrace{t(\mathbf{x})}_{\text{data}} \approx \widehat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$
  - ☐ Why is this called linear regression???
    - ■ model is linear in the parameters

- Precisely, minimize the <span style="color:green">residual squared error</span>:

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

11

---

# The regression problem in matrix notation

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

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



$\mathbf{H} =$ ... $h_1 \ldots h_K$, N data points, K basis functions

$\mathbf{w} =$ ... K basis func, weights

$\mathbf{t} =$ ... $t$, N data points, observations

12

6

# Minimizing the Residual

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

13

---

# Regression solution = simple matrix operations

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

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

$$\text{where } \mathbf{A} = \mathbf{H}^T\mathbf{H} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix} \qquad \mathbf{b} = \mathbf{H}^T\mathbf{t} = \begin{bmatrix} \\ \\ \\ \end{bmatrix}$$

$$\underbrace{\phantom{\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}}}_{\substack{\text{k×k matrix} \\ \text{for k basis functions}}} \qquad \underbrace{\phantom{\begin{bmatrix} \\ \\ \end{bmatrix}}}_{\text{k×1 vector}}$$

14

# But, why?

- Billionaire again, she says: Why sum squared error???
- You say: Gaussians, Gaussians…

- Model: prediction is linear function plus Gaussian noise
  - $t(\mathbf{x}) = \sum_i w_i\, h_i(\mathbf{x}) + \varepsilon_{\mathbf{x}}$

- **Learn $\mathbf{w}$ using MLE**

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

# Maximizing log-likelihood

**Maximize:**

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

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

# Announcements

- Go to recitation!! ☺
  - TBD

- First homework will go out today
  - Due on TBD
  - Start early!!

17

# Bias-Variance Tradeoff

Machine Learning – CSE546

Sham Kakade

University of Washington

Oct 4, 2016

18

# 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

# (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
  - ☐ $Bias^2$ at one point $x$:
  - ☐ Average $Bias^2$:

# 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*:
  - Average variance:

# Bias-Variance Tradeoff

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

# Bias-Variance Decomposition of Error

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

- Expected mean squared error: $\mathrm{MSE} = E_D\left[E_x\left[(t(x) - h_D(x))^2\right]\right]$

- To simplify derivation, drop x:

- Expanding the square:

---

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

- Error can be decomposed:
$$\mathrm{MSE} = E_D\left[E_x\left[(t(x) - h_D(x))^2\right]\right]$$
$$= E_x\left[(t(x) - \bar{h}_N(x))^2\right] + E_D\left[E_x\left[(\bar{h}(x) - h_D(x))^2\right]\right]$$

- 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
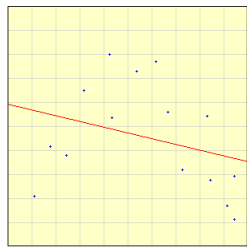- Bias-variance trade-off
- Play with Applet

25

# Overfitting

Machine Learning – CSE546

Sham Kakade

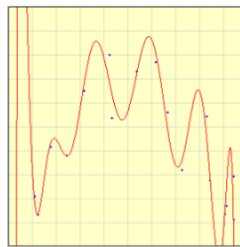University of Washington

Oct 4, 2016

26

# Bias-Variance Tradeoff

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



27

---

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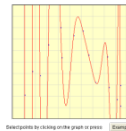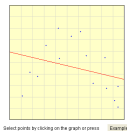- **Training set error:** For a particular set of parameters, loss function on training data:

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

28

14

## Training set error as a function of model complexity

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

29

---

## 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:

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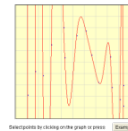\end{aligned}
$$

30

15

## Prediction error as a function of model complexity

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

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

31

---

## Computing prediction error

- Computing prediction
  - Hard integral
  - May not know t(**x**) for every **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 {$\mathbf{x}_1, \ldots, \mathbf{x}_M$} from p(**x**)
  - Approximate integral with sample average

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

32

# 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???

33

---

- 

**Because you cheated!!!**

Training error good estimate for a single **w,**
But you optimized **w** with respect to the training error,
and found **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???

34

17

# 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 – $\{\mathbf{x}_1,\dots,\mathbf{x}_{Ntrain}\}$
  - Test data – $\{\mathbf{x}_1,\dots,\mathbf{x}_{Ntest}\}$
- Use training data to optimize parameters $\mathbf{w}$
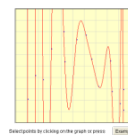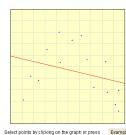- **Test set error:** For the *final output* $\hat{\mathbf{w}}$, evaluate the error using:

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

35

---

# Test set error as a function of model complexity

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

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

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

36

18

# Overfitting

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

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

37

# How many points to I use for training/testing?

- Very hard question to answer!
  - ☐ Too few training points, learned **w** is bad
  - ☐ Too few test points, you never know if you reached a good solution
- Bounds, such as Hoeffding's inequality can help:

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

- 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…
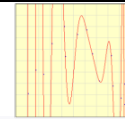    - e.g., bootstrapping

38

19

# Error estimators

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

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

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

39

# 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

40

# Error estimators

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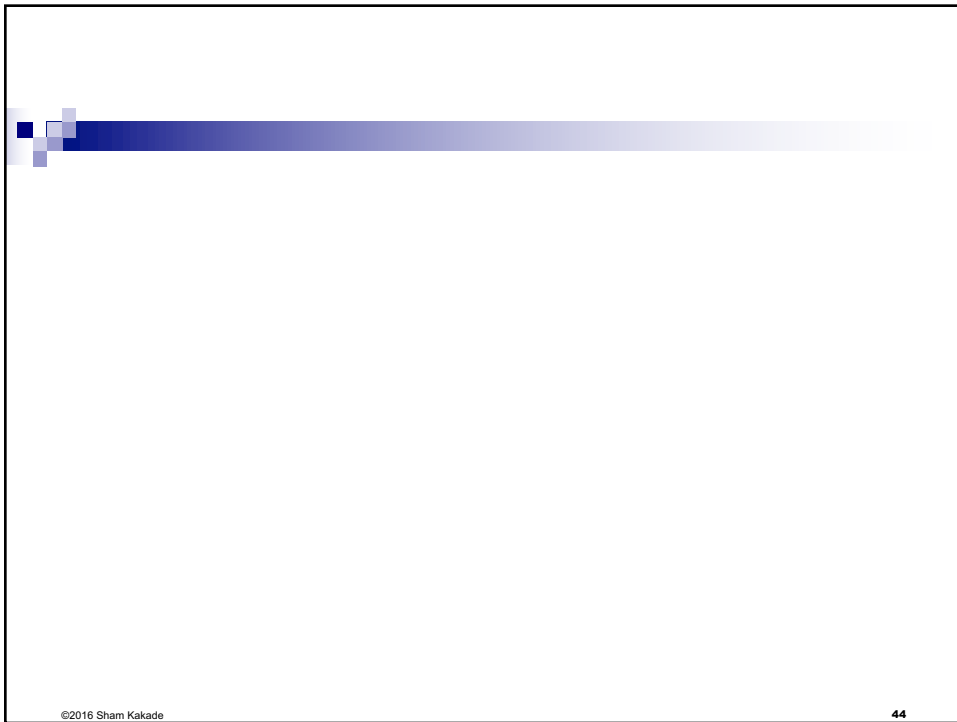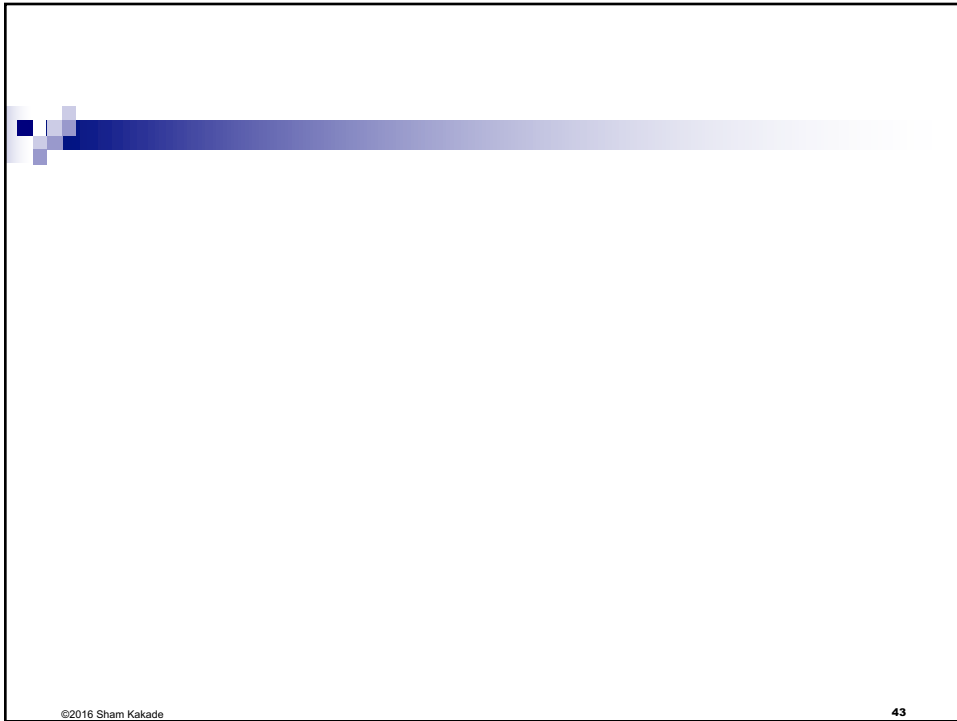do any any any any learning on the test data

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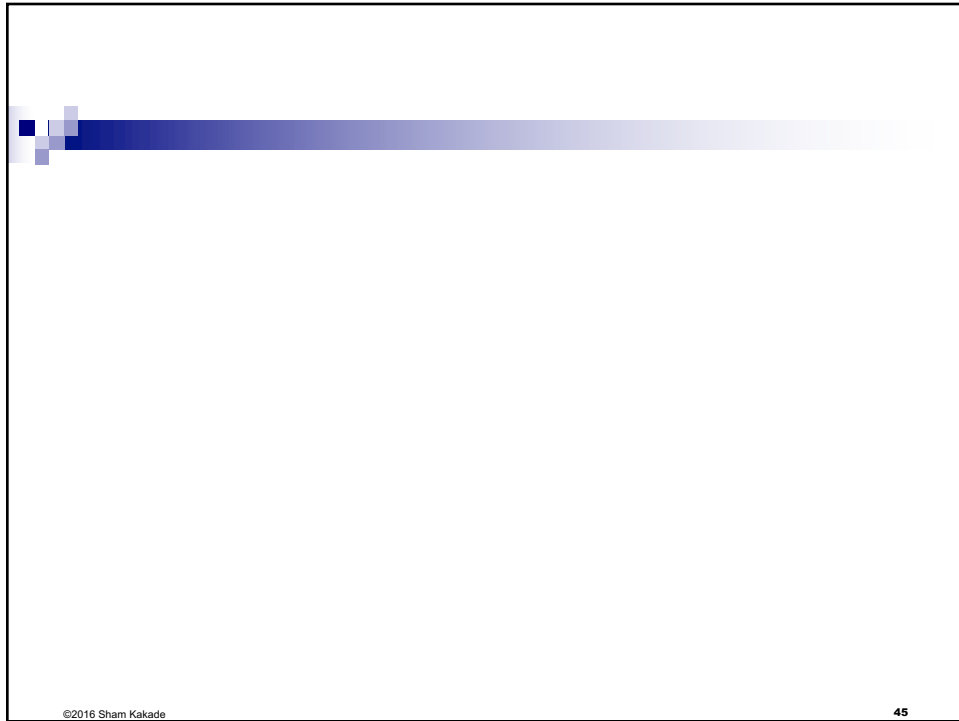(We will address this problem later in the quarter)

$er$

$er$

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

41

---

# 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

42

43

44

22

45

# Bayesian Methods

Machine Learning – CSE546

Sham Kakade

University of Washington

Oct 4, 2016

46

# 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 $\theta$, we obtain a distribution over possible values of $\theta$

47

# Bayesian Learning

- Use Bayes rule:

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

- Or equivalently:

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

48

# 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
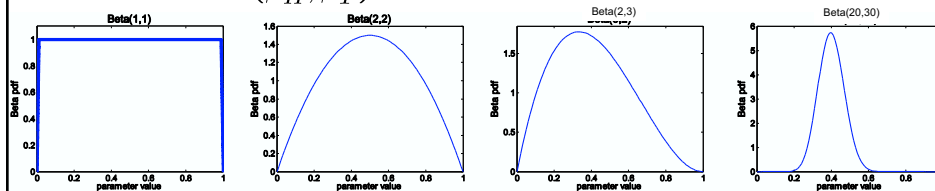  - ☐ **For Binomial, conjugate prior is Beta distribution**

49

---

# Beta prior distribution – P(θ)

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

Mean:

Mode:



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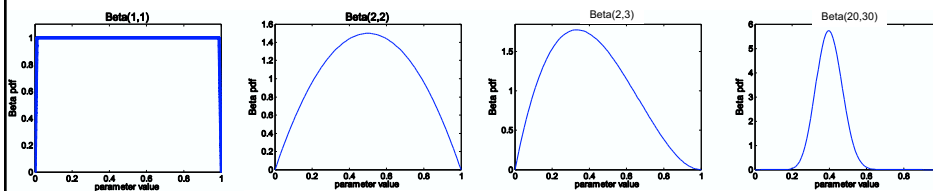- Posterior:  $P(\theta \mid \mathcal{D}) \;\propto\; P(\mathcal{D} \mid \theta)P(\theta)$

50

25

# Posterior distribution

■ Prior: $Beta(\beta_H, \beta_T)$

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

■ Posterior distribution:

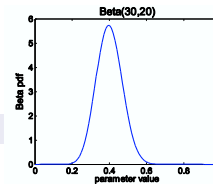$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

---

# Using Bayesian posterior



■ Posterior distribution:

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

■ Bayesian inference:

   ☐ No longer single parameter:

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

   ☐ Integral is often hard to compute
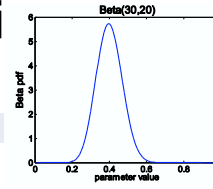
# MAP: Maximum a posteriori approximation

Beta(30,20)

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

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

- As more data is observed, Beta is more certain

- MAP: use most likely parameter:

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

53

---

# MAP for Beta distribution

Beta(30,20)

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

- MAP: use most likely parameter:

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

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

54