# Linear Regression

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

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

---

1

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

3

---

# Learning a Gaussian

HW v. 98
Scores : 96
85
$Y_N$ :

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

- Learn parameters
  - Mean
  - Variance

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} x_i \quad \text{why?} \quad \text{MLE}$$

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

4

2

# MLE for Gaussian

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

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

$$\hat{\mu}_{MLE}, \hat{\sigma}^2_{MLE} = \underset{\mu,\sigma}{argmax}\ P(D|\mu,\sigma) = \underset{\mu,\sigma}{argmax}\ \ln P(D|\mu,\sigma)$$

- Log-likelihood of data:

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

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

5

---

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

- What's MLE for mean?

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

6

# MLE for variance

- Again, set derivative to zero:

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

7

# 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:

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

8

# Prediction of continuous variables

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

9

# The regression problem

- **Instances:** $<\mathbf{x}_j, t_j>$
- **Learn:** Mapping from x to t(**x**)
- **Hypothesis space:**
  - Given, basis functions
  - Find coeffs $\mathbf{w}=\{w_1,\ldots,w_k\}$

$$H = \{h_1, \ldots, h_K\}$$

$$\underbrace{t(\mathbf{x})}_{\text{data}} \approx \widehat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$$

  - Why is this called linear regression???
    - model is linear in the parameters

- Precisely, minimize the residual squared error:

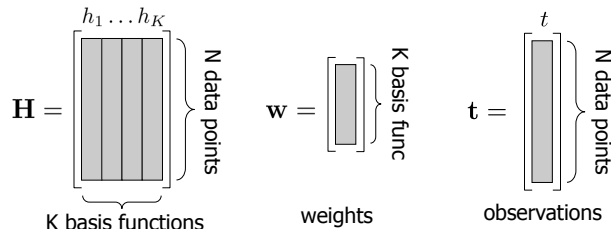$$\mathbf{w}^* = \arg\min_{\mathbf{w}} \sum_j \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$

10

5

## The regression problem in matrix notation

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

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

$\mathbf{H} =$ [matrix] $\quad h_1 \ldots h_K$ / K basis functions / N data points

$\mathbf{w} =$ [vector] K basis func / weights

$\mathbf{t} =$ [vector] $t$ / N data points / observations

©2005-2014 Carlos Guestrin                                       11

## Minimizing the Residual

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

12

6

## Regression solution = simple matrix operations

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

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

$$\text{where } \mathbf{A} = \mathbf{H}^T\mathbf{H} = \underbrace{\begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}}_{\substack{k \times k \text{ matrix} \\ \text{for } k \text{ basis functions}}} \qquad \mathbf{b} = \mathbf{H}^T\mathbf{t} = \underbrace{\begin{bmatrix} \\ \\ \\ \end{bmatrix}}_{k \times 1 \text{ vector}}$$

13

---

## But, why?

- Billionaire (again) says: Why sum squared error???
- You say: Gaussians, Dr. Gateson, Gaussians…

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

- Learn **w** using MLE

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

14

# Maximizing log-likelihood

**Maximize:**

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

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

15

---

# Announcements

- Go to recitation!! ☺
  - □ Wednesday, 5pm in EEB 045

- First homework will go out today
  - □ Due on October 14
  - □ Start early!!

16

# Bias-Variance Tradeoff

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

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

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

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 $x$:
  - ☐ Average variance:

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:  1 ▼   ◉ Fit Y to X
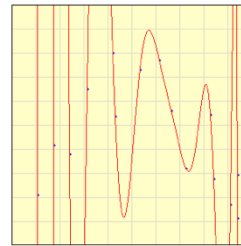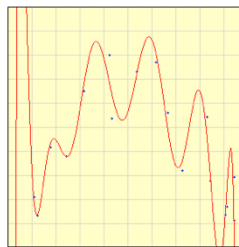                              ○ Fit X to Y

Calculate | View Polynomial | Reset

Select points by clicking on the graph or press    Example

Degree of polynomial:  13 ▼   ◉ Fit Y to X
                               ○ Fit X to Y

Calculate | View Polynomial | Reset

Select points by clicking on the graph or press    Example

Degree of polynomial:  13 ▼   ◉ Fit Y to X
                               ○ Fit X to Y

Calculate | View Polynomial | Reset

©2005-2014 Carlos Guestrin

21

---

# 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:

©2005-2014 Carlos Guestrin

22

11

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

- Error can be decomposed:

$$\text{MSE} = E_D \left[ E_x \left[ (t(x) - h_D(x))^2 \right] \right]$$
$$= E_x \left[ (t(x) - \bar{h}_N(x))^2 \right] + E_D \left[ E_x \left[ (\bar{h}(x) - h_D(x))^2 \right] \right]$$

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

23

---

# 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

24

# Overfitting

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

25

---

# Bias-Variance Tradeoff

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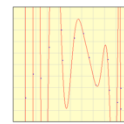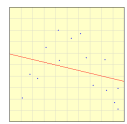  - More complex class → more variance

26

# 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
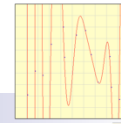- **Training set error:** For a particular set of parameters, loss function on training data:

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

27

---

# Training set error as a function of model complexity

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

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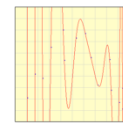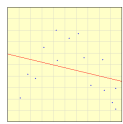
$$
\begin{aligned}
error_{true}(\mathbf{w}) &= E_{\mathbf{x}}\left[\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x})\right)^2\right] \\
&= \int_{\mathbf{x}}\left(t(\mathbf{x}) - \sum_i w_i h_i(\mathbf{x})\right)^2 p(\mathbf{x})d\mathbf{x}
\end{aligned}
$$

29

---

# Prediction error as a function of model complexity

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

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

30

15

# 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 {**x**$_1$,...,**x**$_M$} from p(**x**)
  - Approximate integral with sample average

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

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???

32

# Why training set error doesn't approximate prediction error?

**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???

---

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

$$error_{test}(\mathbf{w}) = \frac{1}{N_{test}} \sum_{j=1}^{N_{test}} \left( t(\mathbf{x}_j) - \sum_i w_i h_i(\mathbf{x}_j) \right)^2$$
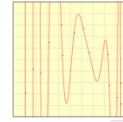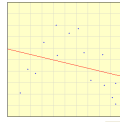
# Test set error as a function of model complexity

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

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

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

35

# Overfitting

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

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

36

# 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…
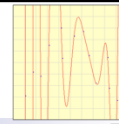    - e.g., bootstrapping

37

# Error estimators

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

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

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

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

39

---

# Error estimators
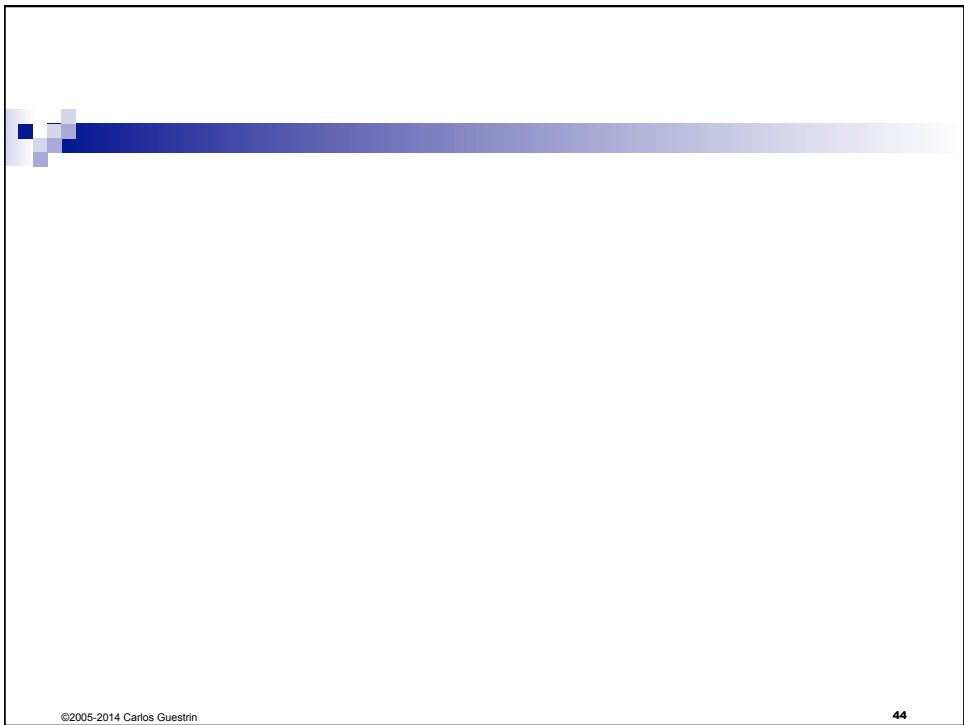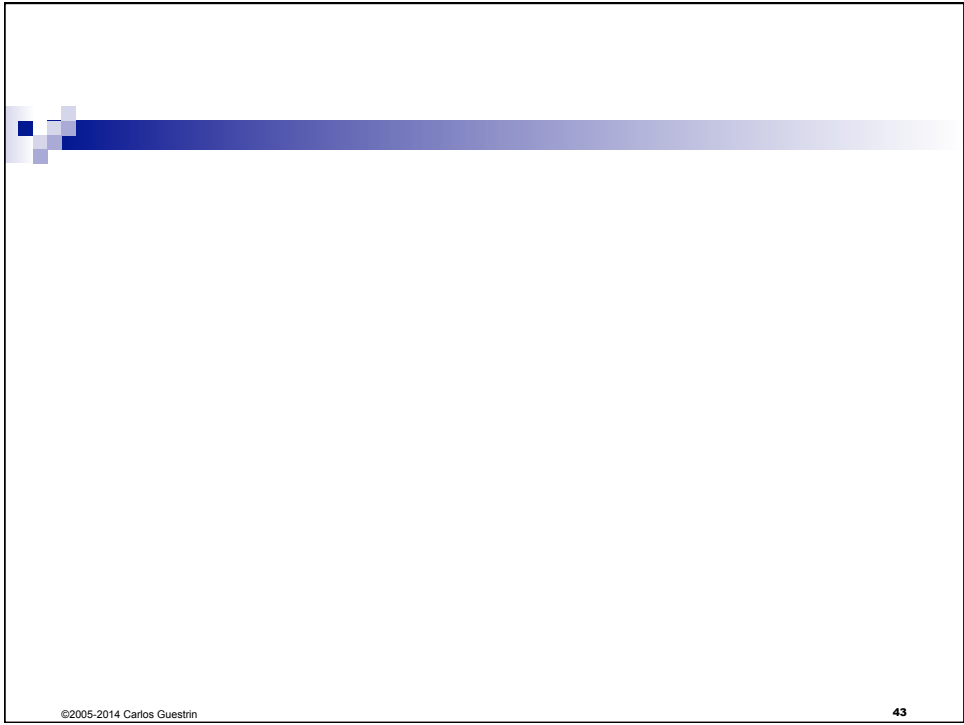
$err$

$err$

**Be careful!!!**

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)

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

40

20

# 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

41

---

42

43

44

# Bayesian Methods

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

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

46

# Bayesian Learning

- Use Bayes rule:

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

- Or equivalently:

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

# 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
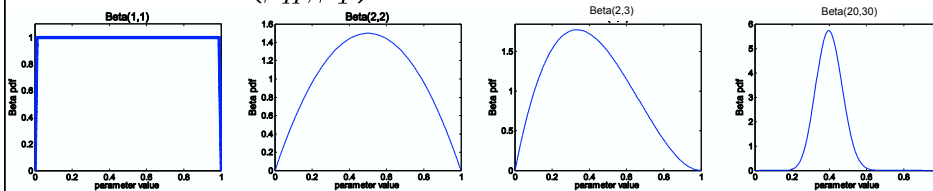  - ☐ Simple posterior form
- Conjugate priors:
  - ☐ Closed-form representation of posterior
  - ☐ **For Binomial, conjugate prior is Beta distribution**

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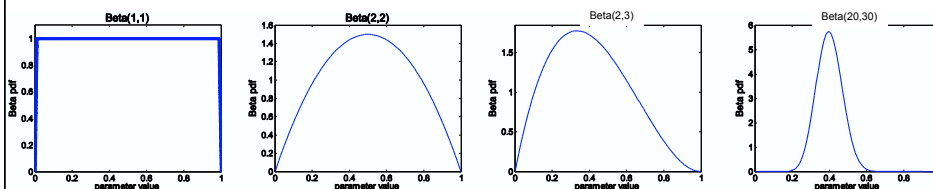- Posterior: $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

49

---

# Posterior distribution

- Prior: $Beta(\beta_H, \beta_T)$
- Data: $\alpha_H$ heads and $\alpha_T$ tails

- Posterior distribution:

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

50

# Using Bayesian posterior

Beta(30,20)

- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$
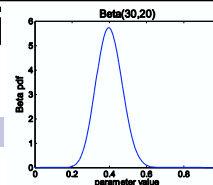
- Bayesian inference:
  - □ No longer single parameter:

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

  - □ Integral is often hard to compute

51

---
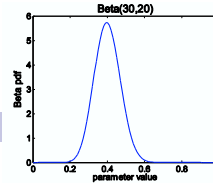
# MAP: Maximum a posteriori approximation

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:

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

52

# MAP for Beta distribution



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

- 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!**

53