

# Bayesian Methods

Machine Learning – CSE546

Carlos Guestrin

University of Washington

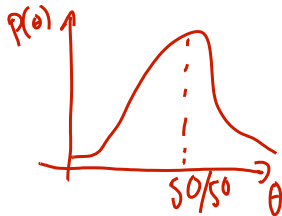
September 30, 2013

©2005-2013 Carlos Guestrin

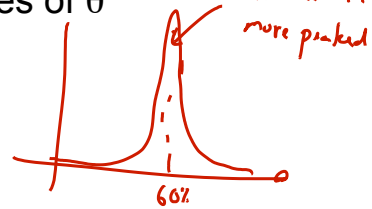
1

## 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  $\theta$



flip coin a few times



©2005-2013 Carlos Guestrin

2

# Bayesian Learning

- Use Bayes rule:

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

*data likelihood*  $\mathcal{D} = (3H, 2T)$   
*prior*  
*posterior*  
*normalizer*

- Or equivalently:

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

*unnormalized version*

MLE:  
 $\max_{\theta} P(\mathcal{D} | \theta)$   
 $\Rightarrow \hat{\theta}_{MLE}$

Bayesian: distribution  $P(\theta | \mathcal{D})$   
 when  $P(\theta)$  is uniform  
 $P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)$

# Bayesian Learning for Thumbtack

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

*beta* *binomial* *beta*

- 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
- philosophy*

- Conjugate priors:

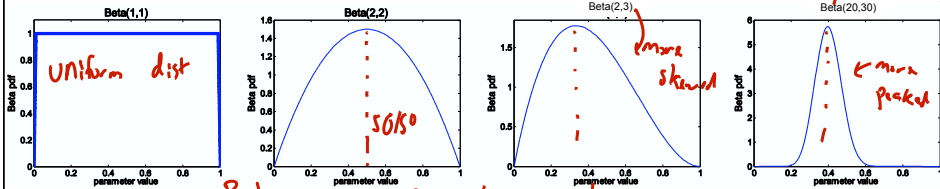
- Closed-form representation of posterior
- For Binomial, conjugate prior is Beta distribution**

$B(\cdot, \cdot)$  ← is the beta function  $\equiv$  normalizer

# 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:  $\beta_H / (\beta_H + \beta_T)$   
 Mode:  $(\beta_H - 1) / (\beta_H + \beta_T - 2)$



Beta is conjugate to Binomial

- Likelihood function:  $P(D | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$
  - Posterior:  $P(\theta | D) \propto P(D | \theta) P(\theta)$  ←  $Beta(\beta_H, \beta_T)$
- $$P(\theta | D = (\alpha_H, \alpha_T)) \propto P(D | \theta) P(\theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T} \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)}$$
- ↳  $\propto \theta^{\alpha_H + \beta_H - 1} (1 - \theta)^{\alpha_T + \beta_T - 1} = \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T)$
- Handwritten note: "doesn't depend on  $\theta$ "*

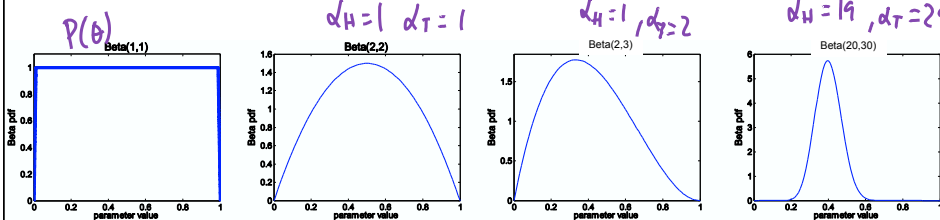
©2005-2013 Carlos Guestrin

5

# Posterior distribution

- Prior:  $\text{Beta}(\beta_H, \beta_T)$
- Data:  $\alpha_H$  heads and  $\alpha_T$  tails
- Posterior distribution:

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

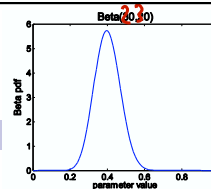


Bayesian updates

©2005-2013 Carlos Guestrin

6

# 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:

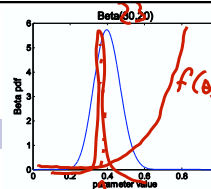
$$E[f(\theta)] = \int_0^1 f(\theta) P(\theta | \mathcal{D}) d\theta$$

integrals can be hard to compute

average over possible parameters

- Integral is often hard to compute

# MAP: Maximum a posteriori approximation



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

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

MAP can be arbitrarily bad approx.

- As more data is observed, Beta is more certain

- MAP: use most likely parameter:

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

possible f:

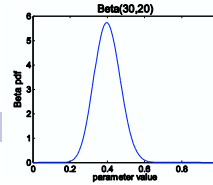
simplest

$$f(\theta) = \theta$$

$$\Rightarrow E[f(\theta)] = E[\theta] = \text{mean } \theta$$

$$f(\theta) = \text{returns } 0 \text{ or } 1 \text{ if true then } \theta$$

# 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:

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

$$\hat{\theta}_{\text{MLE}} = \alpha_H / (\alpha_H + \alpha_T)$$

$(\beta_H + 1, \beta_T + 1) \leftarrow$  take coin flips over MLE

$$\frac{\beta_H + \alpha_H - 1}{\beta_H + \alpha_H + \beta_T + \alpha_T - 2}$$

- Beta prior equivalent to extra thumbtack flips
- As  $N \rightarrow \infty$ , prior is "forgotten"  $\alpha_H, \alpha_T$  dominates  $\beta_H, \beta_T$
- But, for small sample size, prior is important!**

©2005-2013 Carlos Guestrin

9

# Linear Regression

Machine Learning – CSE546

Carlos Guestrin

University of Washington

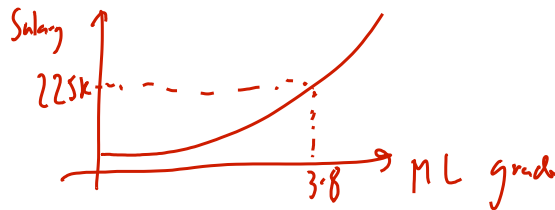
September 30, 2013

©2005-2013 Carlos Guestrin

10

# 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. *and ML grade*
- You sayz: **I can regress that...**



©2005-2013 Carlos Guestrin

11

# The regression problem $\langle 3.9, 225k \rangle$

- **Instances:**  $\langle \mathbf{x}_j, t_j \rangle$
- **Learn:** Mapping from  $\mathbf{x}$  to  $t(\mathbf{x})$   $\hat{f}(\mathbf{x}) \rightarrow t(\mathbf{x})$
- **Hypothesis space:**
  - Given, basis functions,  $H = \{h_1, \dots, h_K\}$  *features*
  - Find coeffs  $\mathbf{w} = \{w_1, \dots, w_K\}$   $t(\mathbf{x}) \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$  *data*  $\hat{f}(\mathbf{x})$  *non-linear*
  - Why is this called linear regression???
  - model is linear in the parameters

*loss function / the risk*

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

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

*target* *estimate* *squared*

©2005-2013 Carlos Guestrin

12

## The regression problem in matrix notation

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

$\downarrow$  same

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

*how do we min?*

$\mathbf{H} =$

K basis functions

$\mathbf{w} =$

weights

$\mathbf{t} =$

observations

©2005-2013 Carlos Guestrin 13

## Minimizing the Residual

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

$\nabla F(\mathbf{w}) = 0 \Rightarrow \nabla H^T (\mathbf{H}\mathbf{w} - \mathbf{t}) = 0$

$H^T \mathbf{H}\mathbf{w} - H^T \mathbf{t} = 0$

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

*ged!!!*

*in scalar calculus*

$$\frac{\partial}{\partial w} (\alpha w - t) (\alpha w - t)$$

$$= 2\alpha (\alpha w - t)$$


---

*in matrix*

$$\frac{\partial}{\partial \mathbf{w}} [(\mathbf{H}\mathbf{w} - \mathbf{t})^T (\mathbf{H}\mathbf{w} - \mathbf{t})]$$

$$= 2 H^T (\mathbf{H}\mathbf{w} - \mathbf{t})$$

©2005-2013 Carlos Guestrin 14

## Regression solution = simple matrix operations

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

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

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

$N < K$   
 $\mathbf{H}^T \mathbf{H}$  is non invertible  
 $k \times k$  matrix  
 for  $k$  basis functions  
 $k \times 1$  vector

©2005-2013 Carlos Guestrin

15

## But, why?

$N(\mu, \sigma^2) \leftarrow$  Gaussian mean  $\mu$   
 variance  $\sigma^2$

- Billionaire (again) says: Why sum squared error???
- You say: Gaussians, Dr. Gateson, Gaussians...
- Model: prediction is linear function plus Gaussian noise

$$\square \mathbf{t}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}) + \epsilon_x$$

$\epsilon_x \sim \text{iid } N(0, \sigma^2)$   
 $\mathbf{t}(\mathbf{x}) \sim \text{iid } N(\sum_i w_i h_i(\mathbf{x}), \sigma^2)$   
 $\epsilon \sim N(0, \sigma^2)$  noise

- Learn  $\mathbf{w}$  using MLE

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

*target*  $\leftarrow$  *input*  $\leftarrow$  *params* *noise variance*

©2005-2013 Carlos Guestrin

16



# Maximizing log-likelihood

$$\ln \Pi = \sum \ln$$

$$\ln e^{\text{smth}} = \text{smth}$$

**Maximize:**

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

*argmax w*

$$= \ln \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N - \sum_{j=1}^N \frac{[t(x_j) - \sum_i w_i h_i(x_j)]^2}{2\sigma^2}$$

*const wrt w*

$$= \underset{w}{\text{argmax}} - \sum_{j=1}^N \frac{[t(x_j) - \sum_{i=1}^k w_i h_i(x_j)]^2}{2\sigma^2}$$

*doesn't affect argmax*

$$= \underset{w}{\text{argmin}} \sum_{j=1}^N [t(x_j) - \sum_{i=1}^k w_i h_i(x_j)]^2$$

*argmax - f = argmin f*

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

# Announcements

- Go to recitation!! ☺
  - Tuesday, 5:30pm in LOW 101
- First homework will go out today
  - Due on October 14
  - Start early!!

# Bias-Variance Tradeoff

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2013

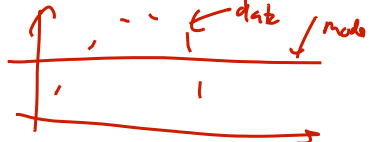
©2005-2013 Carlos Guestrin

19

## 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



©2005-2013 Carlos Guestrin

20

## (Squared) Bias of learner

- Given dataset  $D$  with  $N$  samples, *from  $D \rightarrow$  learn  $h_D(x)$*   
learn function  $h_D(x)$
- If you sample a different dataset  $D'$  with  $N$  samples,  
you will learn different  $h_{D'}(x)$
- **Expected hypothesis:**  $E_D[h_D(x)] = \bar{h}_N(x)$   
*↑ what I expect to learn*
- **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<sup>2</sup> at one point  $x$ :  $(t(x) - \bar{h}_N(x))^2$
  - Average Bias<sup>2</sup>:  $E_x[(t(x) - \bar{h}_N(x))^2]$

©2005-2013 Carlos Guestrin

21

## Variance of learner

$$\text{Var}(x) = E_x[(x - \mu)^2]$$

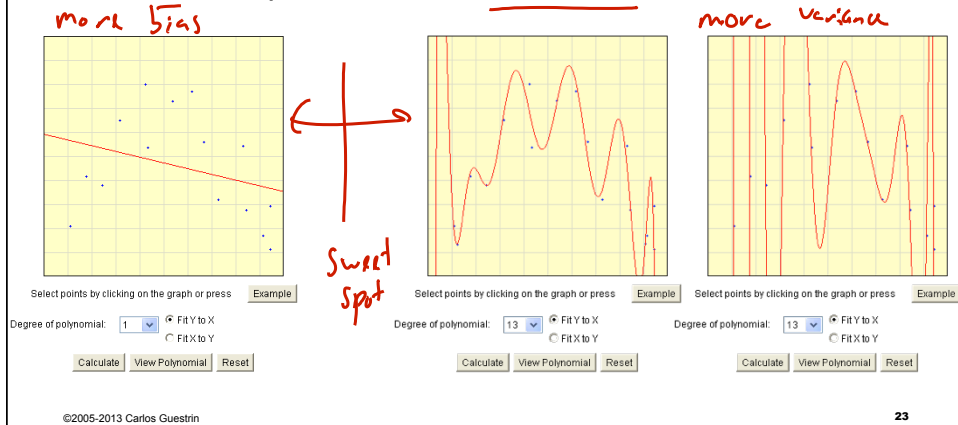
- Given dataset  $D$  with  $N$  samples,  *$D \rightarrow h_D(x)$*   
learn function  $h_D(x)$
- If you sample a different dataset  $D'$  with  $N$  samples,  
you will learn different  $h_{D'}(x)$
- **Variance:** difference between what you expect to learn and  
what you learn from a particular dataset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

©2005-2013 Carlos Guestrin

22

# 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:  $MSE = E_D [E_x [(t(x) - h_D(x))^2]]$

- To simplify derivation, drop x:  $E_D [(t - h_D)^2]$

- Expanding the square: adding & subtracting  $\bar{h}_N$

$$E_D [(t - \bar{h}_N + \bar{h}_N - h_D)^2]$$

$$= E_D [(t - \bar{h}_N)^2] + E_D [(\bar{h}_N - h_D)^2] + 2 E_D [(t - \bar{h}_N)(\bar{h}_N - h_D)]$$

bias
variance

play with this, hint:  
 $h_N = E_D[h_D]$   
 all others are constants

