

# Linear Regression

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

©2005-2014 Carlos Guestrin

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...**

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



©2005-2014 Carlos Guestrin

2

# 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

- $X \sim N(\mu_X, \sigma_X^2)$
- $Y \sim N(\mu_Y, \sigma_Y^2)$
- $Z = X + Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$

©2005-2014 Carlos Guestrin

3

# Learning a Gaussian

- Collect a bunch of data

- Hopefully, i.i.d. samples
- e.g., exam scores

- Learn parameters

- Mean
- Variance

HW scores:  $\begin{bmatrix} 98 \\ 76 \\ 85 \\ \vdots \end{bmatrix}$

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

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

©2005-2014 Carlos Guestrin

4

# MLE for Gaussian

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

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

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

- Log-likelihood of data:

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

©2005-2014 Carlos Guestrin

5

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

- What's MLE for mean?

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

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

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

©2005-2014 Carlos Guestrin

6

## MLE for variance

- Again, set derivative to zero:

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



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

©2005-2014 Carlos Guestrin

7

## Learning Gaussian parameters

- MLE:

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

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

Now  
you  
can  
justify

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

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

©2005-2014 Carlos Guestrin

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...**



©2005-2014 Carlos Guestrin

9

# 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  $\mathbf{w} = \{w_1, \dots, w_K\}$   $t(\mathbf{x}) \approx \hat{f}(\mathbf{x}) = \sum_{i=1}^K w_i h_i(\mathbf{x})$
- Why is this called linear regression???

- model is linear in the parameters

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

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

©2005-2014 Carlos Guestrin

10

## The regression problem in matrix notation

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

*↓ 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

©2005-2014 Carlos Guestrin 11

## Minimizing the Residual

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

*in scalar calculus:*

$$\frac{d}{d\mathbf{w}} [(\alpha\mathbf{w} - \mathbf{t}) (\alpha\mathbf{w} - \mathbf{t})] = 2\alpha (\alpha\mathbf{w} - \mathbf{t})$$

*in matrix:*

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

$\nabla_{\mathbf{w}} F(\mathbf{w}) = 0 \Rightarrow \cancel{2} \mathbf{H}^T (\mathbf{H}\mathbf{w} - \mathbf{t}) = 0$

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

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

©2005-2014 Carlos Guestrin 12

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

$\underbrace{\hspace{10em}}_{k \times k \text{ matrix for } k \text{ basis functions}} \quad \underbrace{\hspace{10em}}_{k \times 1 \text{ vector}}$

©2005-2014 Carlos Guestrin

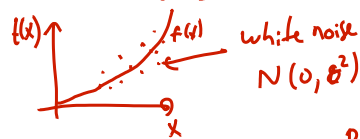
13

## But, why?

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

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

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



- Learn  $\mathbf{w}$  using MLE

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

$\leftarrow$  target  $\quad \leftarrow$  params of mean fn.  $\quad \leftarrow$  noise variance  $\quad \leftarrow$  input

©2005-2014 Carlos Guestrin

14

$\arg \min_w \text{const} + f(w) = \arg \min_w f(w)$        $\ln \pi = \sum \ln$   
 $\ln e^{\text{smth}} = \text{smth}$

## Maximizing log-likelihood

**Maximize:**

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

*arg max w.r.t w*  
*doesn't depend on w*

$$= \arg \max_{\mathbf{w}} \left[ \ln(\cdot) - \sum_{j=1}^N \frac{[t(x_j) - \sum_i w_i h_i(x_j)]^2}{2\sigma^2} \right]$$

$$= \arg \max_{\mathbf{w}} - \sum_{j=1}^N \frac{[t(x_j) - \sum_i w_i h_i(x_j)]^2}{2\sigma^2}$$

$$= \arg \min_{\mathbf{w}} \sum_{j=1}^N \frac{[t(x_j) - \sum_i w_i h_i(x_j)]^2}{2\sigma^2}$$

$\arg \max_{\mathbf{w}} -f(\mathbf{w}) = \arg \min_{\mathbf{w}} f(\mathbf{w})$   
 $\arg \min_{\mathbf{w}} \frac{f(\mathbf{w})}{2\sigma^2} = \arg \min_{\mathbf{w}} f(\mathbf{w})$

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

©2005-2014 Carlos Guestrin 15

## Announcements

- Go to recitation!! 😊
  - Wednesday, 5pm in EEB 045
- First homework will go out today
  - Due on October 14
  - Start early!!

©2005-2014 Carlos Guestrin 16



# Bias-Variance Tradeoff

Machine Learning – CSE546

Carlos Guestrin

University of Washington

September 30, 2014

©2005-2014 Carlos Guestrin

17

## Bias-Variance tradeoff – Intuition

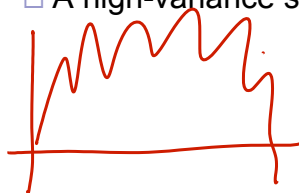
- Model too “simple” → does not fit the data well

- A biased solution



- Model too complex → small changes to the data, solution changes a lot

- A high-variance solution



one new  
data  
point →



©2005-2014 Carlos Guestrin

18

## (Squared) Bias of learner

- Given dataset  $D$  with  $N$  samples, *from  $D \rightarrow$  learn  $h_D$*   
learn function  $h_D(x)$
- If you sample a different dataset  $D'$  with  $N$  samples,  
you will learn different  $h_{D'}(x)$
- **Expected hypothesis:**  $E_D[h_D(x)] = \bar{h}_N$   
*^ 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  $x$ :  $(f(x) - \bar{h}_N(x))^2$
  - Average Bias<sup>2</sup>:  $E_x [(f(x) - \bar{h}_N(x))^2]$

©2005-2014 Carlos Guestrin

19

## Variance of learner

- Given dataset  $D$  with  $N$  samples,  *$D \rightarrow h_D$*   
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  $x$ :  $E_{D'}(h_{D'}(x) - \bar{h}_N(x))^2$   
*^ actually learn*      *^ expect to learn*
  - Average variance:

©2005-2014 Carlos Guestrin

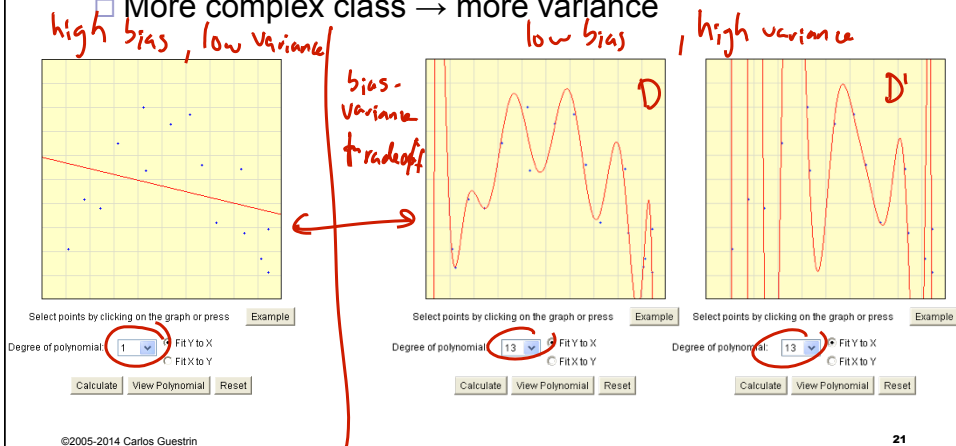
20

# 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

- 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: Add & subtract  $\bar{h}_N$

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

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

$\underbrace{\hspace{10em}}_{\text{bias}^2}$ 
 $\underbrace{\hspace{10em}}_{\text{variance}}$

play with this at home, hint:  
 $\bar{h}_N = E_D [h_D]$   
 all other terms const.

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

- Error can be decomposed:

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

*want to minimize* (pointing to MSE)  
*bias* (under the first term)  
*Variance* (under the second term)

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

©2005-2014 Carlos Guestrin

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

©2005-2014 Carlos Guestrin

24

# Overfitting

Machine Learning – CSE546

Carlos Guestrin

University of Washington

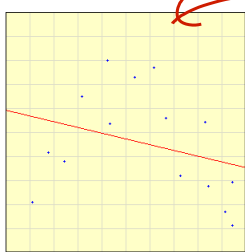
September 30, 2014

©2005-2014 Carlos Guestrin

25

## Bias-Variance Tradeoff

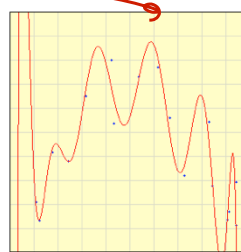
- Choice of hypothesis class introduces learning bias
  - More complex class → less bias
  - More complex class → more variance



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

Degree of polynomial:   Fit Y to X  
 Fit X to Y

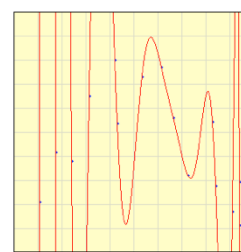
[Calculate](#) [View Polynomial](#) [Reset](#)



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

Degree of polynomial:   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:   Fit Y to X  
 Fit X to Y

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

©2005-2014 Carlos Guestrin

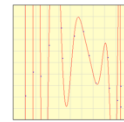
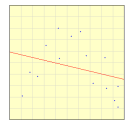
26

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

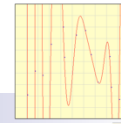
- Given a dataset (Training data)
- Choose a loss function
  - e.g., squared error ( $L_2$ ) for regression
- **Training set error:** For a particular set of parameters, loss function on training data:

$$\text{error}_{train}(\mathbf{w}) = \frac{1}{N_{train}} \sum_{j=1}^{N_{train}} \left( \underset{\substack{\uparrow \\ \text{truth}}}{t(\mathbf{x}_j)} - \sum_i \underset{\substack{\uparrow \\ \text{estimate}}}{w_i h_i(\mathbf{x}_j)} \right)^2$$

# Training set error as a function of model complexity



# 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:

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

# Prediction error as a function of model complexity

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

