

Readings listed in class website

Gaussians

Linear Regression

Bias-Variance Tradeoff

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

September 12th, 2007

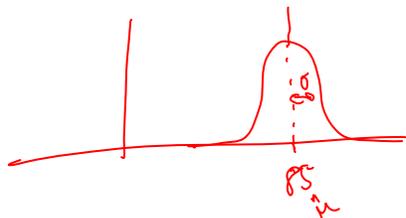
1

©Carlos Guestrin 2005-2007

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...**

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



91
92
89
⋮

2

©Carlos Guestrin 2005-2007

Some properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)

- $X \sim N(\mu, \sigma^2)$

- $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$

$\begin{matrix} - & \uparrow & \uparrow & \\ & 2 & 9 & \\ & & & -9.3 \end{matrix}$

- Sum of Gaussians (independent)

- $X \sim N(\mu_X, \sigma_X^2)$

- $Y \sim N(\mu_Y, \sigma_Y^2)$

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

3

©Carlos Guestrin 2005-2007

Learning a Gaussian

- Collect a bunch of data

- Hopefully, i.i.d. samples

- e.g., exam scores

$x_i \rightarrow \begin{matrix} 91 \\ 87 \\ \vdots \end{matrix}$

- Learn parameters

- Mean = $E[X] = \frac{1}{n} \sum x_i$

- Variance = $E[(X - \mu)^2]$

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

4

©Carlos Guestrin 2005-2007

MLE for Gaussian

$$\ln a^b = b \ln a$$

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

$$D = \{17, 91, 82, \dots\}$$

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

$$P(D | \mu, \sigma) = \prod P(x_i | \mu, \sigma) \quad \text{independence}$$

- Log-likelihood of data:

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

5

©Carlos Guestrin 2005-2007

Your second learning algorithm:

MLE for mean of a Gaussian

$$\frac{d}{dx} (x - \mu)^2 = -2(x - \mu)$$

- What's MLE for mean?

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

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

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

6

©Carlos Guestrin 2005-2007

MLE for variance

- Again, set derivative to zero:

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

$$\begin{aligned} \frac{-N}{\sigma} + \sum_{i=1}^N \frac{2(x_i - \mu)^2}{2\sigma^3} &= 0 \\ \Rightarrow \hat{\sigma}_{MLE}^2 &= \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} \end{aligned}$$

©Carlos Guestrin 2005-2007

7

Learning Gaussian parameters

- MLE:

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

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

- BTW. MLE for the variance of a Gaussian is **biased**

- Expected result of estimation is **not** true parameter!
- Unbiased variance estimator:

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

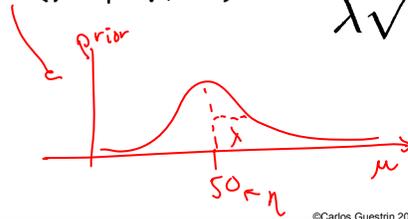
©Carlos Guestrin 2005-2007

Bayesian learning of Gaussian parameters

- Conjugate priors
 - Mean: Gaussian prior
 - Variance: Wishart Distribution

■ Prior for mean:

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



9

MAP for mean of Gaussian

$P(\mu) \propto \frac{1}{\lambda\sqrt{2\pi}}$

$P(\mu | D) \propto P(\mu) \cdot P(D | \mu, \sigma)$

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

$$0 = \frac{d}{d\mu} [\ln P(D | \mu) P(\mu)] = \frac{d}{d\mu} [\ln P(D | \mu) + \ln P(\mu)]$$

ln posterior *likelihood* *prior*

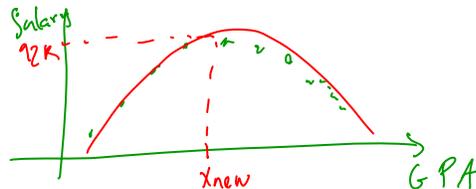
$$\begin{aligned} & \frac{d}{d\mu} \ln P(D | \mu) + \frac{d}{d\mu} \ln P(\mu) \\ &= \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2} + \frac{d}{d\mu} \left[\ln \frac{1}{\lambda\sqrt{2\pi}} - \frac{(\mu - \eta)^2}{2\lambda^2} \right] \\ &= \sum_{i=1}^N \frac{(x_i - \mu)}{\sigma^2} - \frac{\mu - \eta}{\lambda^2} = 0 \end{aligned}$$

$$\begin{aligned} \sum_{i=1}^N \frac{x_i}{\sigma^2} + \frac{\eta}{\lambda^2} &= \frac{N\mu}{\sigma^2} + \frac{\mu}{\lambda^2} \\ \hat{\mu}_{MAP} &= \frac{\sum_{i=1}^N \frac{x_i}{\sigma^2} + \frac{\eta}{\lambda^2}}{\frac{N}{\sigma^2} + \frac{1}{\lambda^2}} \end{aligned}$$

10,

Prediction of continuous variables

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



$\langle 3.8, 95K \rangle$
 $\langle 4.0, 82K \rangle$
 $\langle 3.7, 500K \rangle$
 $\uparrow \quad \quad \uparrow$
 $x \quad \quad y$

11

©Carlos Guestrin 2005-2007

The regression problem

- **Instances:** $\langle x_j, t_j \rangle$ $\leftarrow \langle 3.8, 92K \rangle$
- **Learn:** Mapping from x to $t(x)$ $x \mapsto \mathbb{R}$
- **Hypothesis space:** $H = \{h_1, \dots, h_K\}$
 - Given, basis functions
 - Find coeffs $\mathbf{w} = \{w_1, \dots, w_k\}$ $t(x) \approx \hat{f}(x) = \sum_i w_i h_i(x)$
 - Why is this called linear regression???
 - model is linear in the parameters
- Precisely, minimize the **residual squared error:**

poly:
 $1, x, x^2, x^3, \dots$
 Fourier:
 $\sin(\pi x)$
 $\cos(\pi x)$
 $\sin(2\pi x)$
 \vdots
 Features
 $(t(x_j) - \sum w_i h_i(x_j))$
 Loss function

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

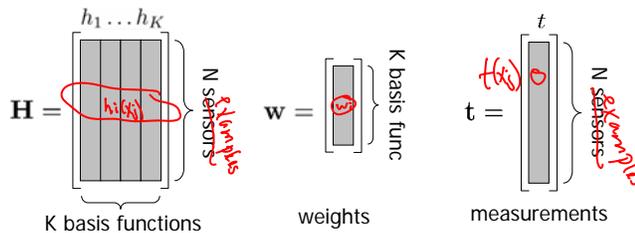
12

©Carlos Guestrin 2005-2007

The regression problem in matrix notation

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

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



13

Regression solution = simple matrix operations

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

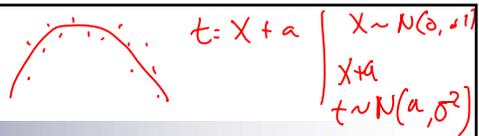
$$\text{solution: } \mathbf{w}^* = \underbrace{(\mathbf{H}^T \mathbf{H})^{-1}}_{\mathbf{A}^{-1}} \underbrace{\mathbf{H}^T \mathbf{t}}_{\mathbf{b}} = \mathbf{A}^{-1} \mathbf{b}$$

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

\mathbf{A} is a $k \times k$ matrix for k basis functions.
 \mathbf{b} is a $k \times 1$ vector.

14

But, why?



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

$$t = \underbrace{\sum_i w_i h_i(\mathbf{x})}_{\mu} + \varepsilon \quad \leftarrow N(0, \sigma^2)$$

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

- Learn \mathbf{w} using MLE

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

15

©Carlos Guestrin 2005-2007

Maximizing log-likelihood

↑ *max*

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

constant doesn't depend on w

$$\ln \left(\frac{1}{\sigma\sqrt{2\pi}} \right)^N \underbrace{- \sum_{j=1}^N [t_j - \sum_i w_i h_i(\mathbf{x}_j)]^2}_{\text{minimize}}$$

iid

$$\equiv \min_{\mathbf{w}} \sum_{j=1}^N (t_j - \sum_i w_i h_i(\mathbf{x}_j))^2$$

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

©Carlos Guestrin 2005-2007

Applications Corner 1

- Predict stock value over time from
 - past values
 - other relevant vars
 - e.g., weather, demands, etc.

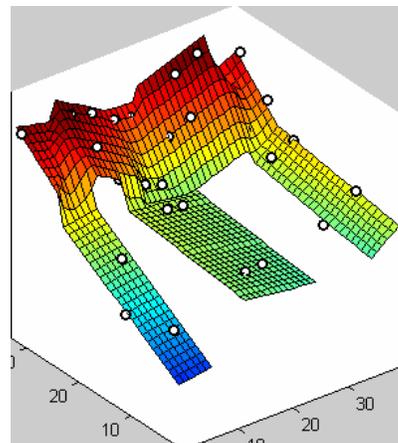
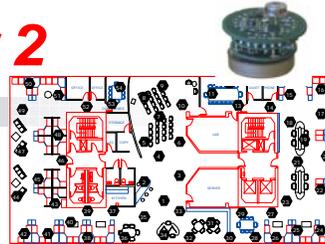


17

©Carlos Guestrin 2005-2007

Applications Corner 2

- Measure temperatures at some locations
- Predict temperatures throughout the environment



[Guestrin et al. '04]

©Carlos Guestrin 2005-2007

Applications Corner 3

■ Predict when a sensor will fail

□ based several variables

- age, chemical exposure, number of hours used,...

→ when will
sensor fail

19

©Carlos Guestrin 2005-2007

Announcements 1

■ Readings associated with each class

- See course website for specific sections, extra links, and further details
- Visit the website frequently

■ Recitations

- Thursdays, 5:00-6:20 in Wean Hall 5409

■ Special recitation on Matlab

- Sept. 18 Tue. 4:30-5:50pm NSH 3002

■ Carlos away on Monday Sept. 17th

- Prof. Eric Xing will teach the lecture

20

©Carlos Guestrin 2005-2007

Announcement 2

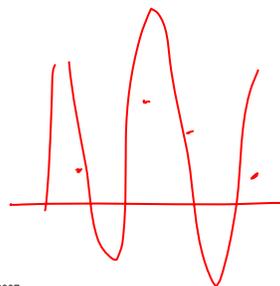
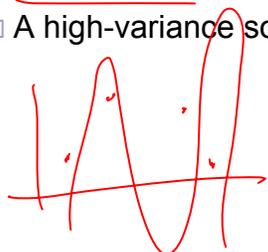
- First homework out later today!
 - Download from course website!
 - Start early!!! :)
 - Due Oct 3rd
- To expedite grading:
 - there are 4 questions
 - please hand in 4 stapled separate parts, one for each question

21

©Carlos Guestrin 2005-2007

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



22

©Carlos Guestrin 2005-2007

(Squared) Bias of learner

- Given dataset D with m samples, learn function $h(x)$
- If you sample a different datasets, m' you will learn different $h(x)$
- **Expected hypothesis:** $E_D[h(x)]$ *on average*
- **Bias:** difference between what you expect to learn and truth
 - Measures how well you expect to represent true solution
 - Decreases with more complex model

$$\text{bias}^2 = \int_x \{ \underbrace{E_D[h(x)]}_{\text{average } h} - \underbrace{t(x)}_{\text{truth}} \}^2 p(x) dx$$

dist. over where you will be tested

23

©Carlos Guestrin 2005-2007

Variance of learner

- Given a dataset D with m samples, you learn function $h(x)$
- If you sample a different datasets, m' you will learn different $h(x)$
- **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

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

$$\text{variance} = \int E_D[(h(x) - \bar{h}(x))^2] p(x) dx$$

specific parameter for D *average parameter*

24

©Carlos Guestrin 2005-2007

Bias-Variance Tradeoff

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

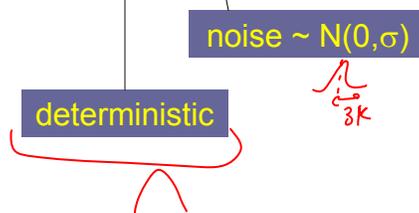
25

©Carlos Guestrin 2005-2007

Bias-Variance decomposition of error

- Consider simple regression problem $f: X \rightarrow T$

$$t = f(x) = g(x) + \varepsilon \quad \approx h(x)$$



Collect some data, and learn a function $h(x)$

What are sources of prediction error? $(f(x) - h(x))^2$

26

©Carlos Guestrin 2005-2007

Sources of error 1 – noise

- What if we have perfect learner, infinite data?
 - If our learning solution $h(x)$ satisfies $h(x)=g(x)$
 - Still have remaining, unavoidable error of σ^2 due to noise ϵ

$$\text{error}(h) = \int_x \int_t (h(x) - t)^2 p(f(x) = t|x)p(x) dt dx$$

27

©Carlos Guestrin 2005-2007

Sources of error 2 – Finite data

- What if we have imperfect learner, or only m training examples?
- What is our expected squared error per example?
 - Expectation taken over random training sets D of size m , drawn from distribution $P(X,T)$

$$E_D \left[\int_x \int_t \{h(x) - t\}^2 p(f(x) = t|x)p(x) dt dx \right]$$

28

©Carlos Guestrin 2005-2007

Bias-Variance Decomposition of Error

Bishop Chapter 3

Assume target function: $t = f(x) = g(x) + \varepsilon$

Then expected sq error over fixed size training sets D drawn from $P(X, T)$ can be expressed as sum of three components:

$$E_{\text{var}}(h) = E_D \left[\int_x \int_t (h(x) - t)^2 p(t|x) p(x) dt dx \right]$$
$$= \text{unavoidableError} + \text{bias}^2 + \text{variance}$$

Where:

$$\text{unavoidableError} = \sigma^2$$

$$\text{bias}^2 = \int (E_D[h(x)] - g(x))^2 p(x) dx$$

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

$$\text{variance} = \int E_D[(h(x) - \bar{h}(x))^2] p(x) dx$$

29

©Carlos Guestrin 2005-2007

What you need to know

- Gaussian estimation
 - MLE
 - Bayesian learning
 - MAP
 - Regression
 - Basis function = features
 - Optimizing sum squared error
 - Relationship between regression and Gaussians
 - Bias-Variance trade-off
 - Play with Applet
- change amount of data
- change degree

30

©Carlos Guestrin 2005-2007