

Point Estimation

Machine Learning – CSE446

Carlos Guestrin

University of Washington

April 3, 2013

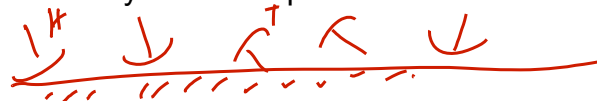
©2005-2013 Carlos Guestrin

1

Your first consulting job

- A billionaire from the suburbs of Seattle asks you a question:

- ☐ He says: I have thumbtack, if I flip it, what's the probability it will fall with the nail up?
- ☐ You say: Please flip it a few times:



$$p(H) = \frac{3}{5}$$

- ☐ You say: The probability is:
- ☐ **He says: Why???**
- ☐ You say: Because...

©2005-2013 Carlos Guestrin

2

Thumbtack – Binomial Distribution

- $P(\text{Heads}) = \theta$, $P(\text{Tails}) = 1 - \theta$

$$P(\mathcal{D} | \theta) = P(\text{HHTTH}) = \theta \theta (1 - \theta) (1 - \theta) \theta = \theta^3 (1 - \theta)^2$$

- Flips are i.i.d.:
 - Independent events
 - Identically distributed according to Binomial distribution
- Sequence \mathcal{D} of α_H Heads and α_T Tails

$$P(\mathcal{D} | \theta) = \theta^{\alpha_H} (1 - \theta)^{\alpha_T}$$

model
observed
HHTTH
learning task
choose $\hat{\theta}$
pick θ that
maximizes the
likelihood of
seeing these
observations
according to
my model
MLE

©2005-2013 Carlos Guestrin

3

Maximum Likelihood Estimation

- **Data:** Observed set \mathcal{D} of α_H Heads and α_T Tails
- **Hypothesis:** Binomial distribution
- Learning θ is an optimization problem
 - What's the objective function?
- MLE: Choose θ that maximizes the probability of observed data:

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \end{aligned}$$

©2005-2013 Carlos Guestrin

4

Your first learning algorithm

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \ln P(\mathcal{D} | \theta) \\ &= \arg \max_{\theta} \ln \theta^{\alpha_H} (1 - \theta)^{\alpha_T}\end{aligned}$$

- Set derivative to zero: $\frac{d}{d\theta} \ln P(\mathcal{D} | \theta) = 0$

©2005-2013 Carlos Guestrin

5

How many flips do I need?

$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

- Billionaire says: I flipped 3 heads and 2 tails.
- You say: $\theta = 3/5$, I can prove it!
- He says: What if I flipped 30 heads and 20 tails?
- You say: Same answer, I can prove it!
- **He says: What's better?**
- You say: Humm... The more the merrier???
- He says: Is this why I am paying you the big bucks???

©2005-2013 Carlos Guestrin

6

Simple bound (based on Hoeffding's inequality)

- For $N = \alpha_H + \alpha_T$, and $\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$

- Let θ^* be the true parameter, for any $\epsilon > 0$:

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

PAC Learning

- PAC: Probably Approximate Correct
- Billionaire says: I want to know the thumbtack parameter θ , within $\epsilon = 0.1$, with probability at least $1 - \delta = 0.95$. How many flips?

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

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...**

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

©2005-2013 Carlos Guestrin

9

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

©2005-2013 Carlos Guestrin

10

Learning a Gaussian

- Collect a bunch of data

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

- Learn parameters

- Mean
- Variance

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

©2005-2013 Carlos Guestrin

11

MLE for Gaussian

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

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

- Log-likelihood of data:

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

©2005-2013 Carlos Guestrin

12

Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

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

MLE for variance

- Again, set derivative to zero:

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

Learning Gaussian parameters

- MLE:

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

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

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

- ☐ Expected result of estimation is **not** true parameter!

- ☐ Unbiased variance estimator:

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

©2005-2013 Carlos Guestrin

15

What you need to know...

- Learning is...

- ☐ Collect some data
 - E.g., thumbtack flips
- ☐ Choose a hypothesis class or model
 - E.g., binomial
- ☐ Choose a loss function
 - E.g., data likelihood
- ☐ Choose an optimization procedure
 - E.g., set derivative to zero to obtain MLE
- ☐ Collect the big bucks

- Like everything in life, there is a lot more to learn...

- ☐ Many more facets... Many more nuances...
- ☐ The fun will continue...

©2005-2013 Carlos Guestrin

16

Announcement: R Tutorial

- R: open-source scripting language for stats & ML
- A lot of resources online
- Tutorial:
 - When: Thursday April 4th at 6:00pm
 - Where: EEB 125
- Before attending please download and install R:
 - <http://www.r-project.org/>
- We recommend using an R environment such as:
 - R studio: <http://www.rstudio.com/>
 - Tinn-R: <http://www.sciviews.org/Tinn-R/>

©2005-2013 Carlos Guestrin

17

Linear Regression Bias-Variance Tradeoff

Machine Learning – CSE446
Carlos Guestrin
University of Washington
April 3, 2013

©2005-2013 Carlos Guestrin

18

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...**

©2005-2013 Carlos Guestrin

19

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 $\mathbf{w} = \{w_1, \dots, w_K\}$ $\underbrace{t(\mathbf{x})}_{\text{data}} \approx \hat{f}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x})$
 - Why is this called linear regression???
 - model is linear in the parameters
- Precisely, minimize the **residual squared error**:

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

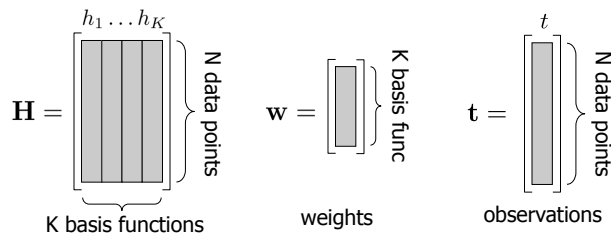
©2005-2013 Carlos Guestrin

20

The regression problem in matrix notation

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

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



21

Minimizing the Residual

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

©2005-2013 Carlos Guestrin

22

Regression solution = simple matrix operations

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

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

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

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

©2005-2013 Carlos Guestrin

23

But, why?

- Billionaire (again) says: Why sum squared error???
- You say: Gaussians, Dr. Gateson, Gaussians...
- Model: prediction is linear function plus Gaussian noise
 - $\mathbf{t}(\mathbf{x}) = \sum_i w_i h_i(\mathbf{x}) + \epsilon_{\mathbf{x}}$

- Learn \mathbf{w} using MLE

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

©2005-2013 Carlos Guestrin

24

Maximizing log-likelihood

Maximize:

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

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

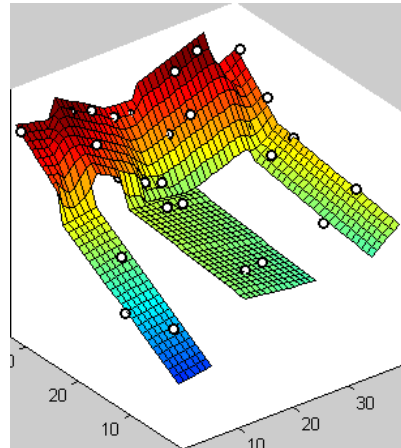
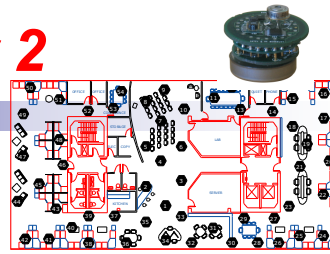
Applications Corner 1

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



Applications Corner 2

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



[Guestrin et al. '04]

©2005-2013 Carlos Guestrin

27

Applications Corner 3

- Predict when a sensor will fail
 - based several variables
 - age, chemical exposure, number of hours used,...

©2005-2013 Carlos Guestrin

28

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

29

(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
 - Bias² at one point x :
 - Average Bias²:

©2005-2013 Carlos Guestrin

30

Variance of learner

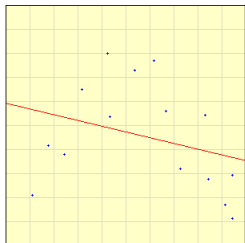
- Given dataset D with N samples, learn function $h_D(x)$
- If you sample a different dataset D' with N samples, you will learn different $h_{D'}(x)$
- **Variance:** difference between what you expect to learn and what you learn from a particular dataset
 - Measures how sensitive learner is to specific dataset
 - Decreases with simpler model
 - Variance at one point x :
 - Average variance:

©2005-2013 Carlos Guestrin

31

Bias-Variance Tradeoff

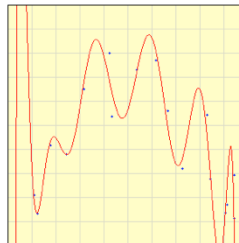
- Choice of hypothesis class introduces learning bias
 - More complex class \rightarrow less bias
 - More complex class \rightarrow more variance



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

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

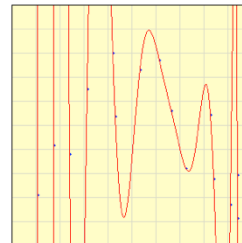
[Calculate](#) [View Polynomial](#) [Reset](#)



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](#)

©2005-2013 Carlos Guestrin

32

Bias-Variance Decomposition of Error

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

- Expected mean squared error: $\text{MSE} = E_D \left[E_x \left[(t(x) - h_D(x))^2 \right] \right]$
- To simplify derivation, drop x :
- Expanding the square:

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

- Error can be decomposed:
$$\begin{aligned}\text{MSE} &= E_D \left[E_x \left[(t(x) - h_D(x))^2 \right] \right] \\ &= E_x \left[(t(x) - \bar{h}_N(x))^2 \right] + E_D \left[E_x \left[(\bar{h}_N(x) - h_D(x))^2 \right] \right]\end{aligned}$$
- Choice of hypothesis class introduces learning bias
 - More complex class \rightarrow less bias
 - More complex class \rightarrow 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