

Gaussians and Linear Regression

Machine Learning - 10601

Geoff Gordon, Miroslav Dudík

<http://www.cs.cmu.edu/~ggordon/10601/>

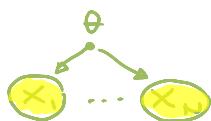
September 16, 2009

Last time: density estimation

Input: X_1, X_2, \dots, X_N i.i.d.

Output: $p(X)$

Maximum likelihood:



$$\max_{\theta} p(\text{data} | \theta) = \prod_{n=1}^N p(x_n | \theta)$$

$$\hat{\theta}_{ML} \leftarrow \max_{\theta} \sum_n \log p(x_n | \theta)$$

GOODNESS OF FIT

Maximum a posteriori (MAP): $\max_{\theta} p(\theta | \text{data})$



$$\hat{\theta}_{MAP} \leftarrow \max_{\theta} \sum_n \log p(x_n | \theta) + \log p(\theta)$$

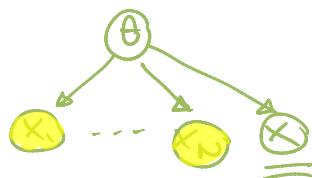
GOODNESS OF FIT + LEVEL OF BELIEF

Last time: density estimation

Input: X_1, X_2, \dots, X_N

Output: $p(X)$

Bayesian prediction:



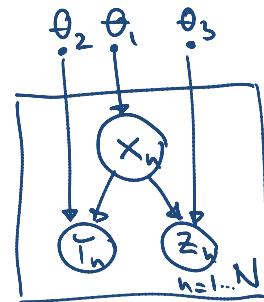
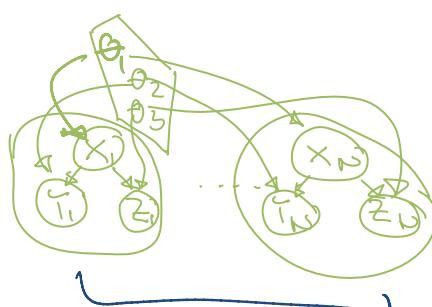
$$p(x | \underbrace{x_1, \dots, x_n}_{\text{data}})$$

$$\propto \int p(x|\theta) p(\theta | \text{data}) d\theta$$

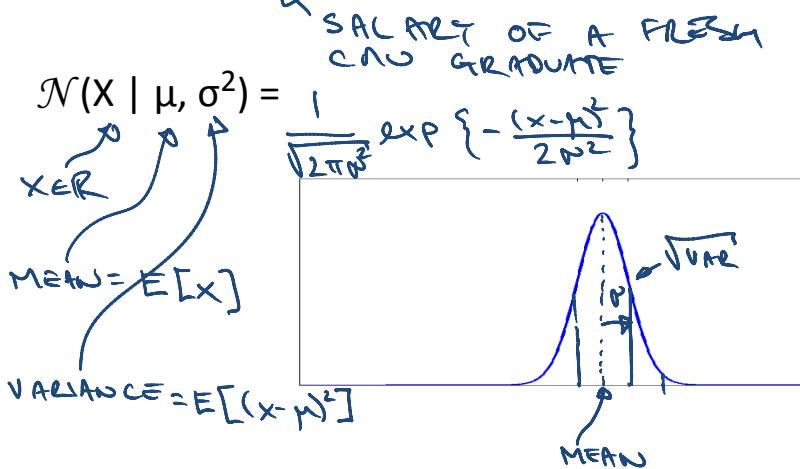
IN GENERAL: DIFFICULT
CONJUGATE PRIORS
BINOM & **BETA**

Density estimation

- we used the thumbtack example
- similar for structured variables
 - diseases and symptoms; natural images; natural language...



Continuous variables: Gaussian



Obtaining Gaussians from Gaussians

- adding constants and scaling

$$X \sim \mathcal{N}(X | \mu, \sigma^2)$$

$$Y = aX + b$$

$$Y \sim \mathcal{N}(Y | a\mu + b, a^2\sigma^2)$$

↑
NORMAL DIST.

- summing independent Gaussians

$$X \sim \mathcal{N}(X | \underline{\mu_x}, \sigma_x^2)$$

$$Y \sim \mathcal{N}(Y | \underline{\mu_y}, \sigma_y^2)$$

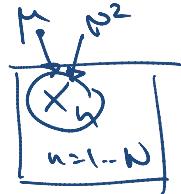
$$Z = X + Y$$

$$Z \sim \mathcal{N}(Z | \underline{\mu_x + \mu_y}, \sigma_x^2 + \sigma_y^2)$$

$$E[Z] = E[X] + E[Y]$$

Learning Gaussians

Maximum Likelihood



$$\max_{\mu, \sigma^2} \sum_n \log p(x_n | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$\sum_n \left[-\frac{1}{2} \log(2\pi\sigma^2) - \frac{(x_n - \mu)^2}{2\sigma^2} \right]$$

CONCAVE

$$\frac{\partial}{\partial \mu} = \sum_n \left[+ \frac{2(x_n - \mu)}{2\sigma^2} \right] = 0$$

$$E[\hat{\mu}_{ML}] = \mu$$

$$\sum_n x_n = N\mu$$

$$\frac{\partial}{\partial \sigma^2} = \rightarrow \hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_n (x_n - \hat{\mu}_{ML})^2$$

BIASED ESTIMATE

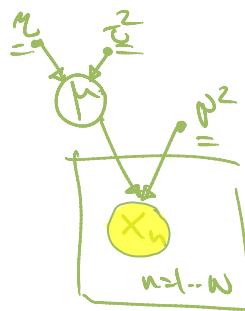
$$E[\hat{\sigma}_{ML}^2] = \frac{N-1}{N} \sigma^2$$

Posterior of the mean

$$x \sim \mathcal{N}(x | \mu, \sigma^2)$$

$$\mu \sim \mathcal{N}(\mu | \gamma, \tau^2)$$

KNOWN



$$p(\mu | x_1, \dots, x_n, \gamma, \tau^2, \sigma^2)$$

$$= p(\mu | \gamma, \tau^2) \underbrace{p(x_1, \dots, x_n | \mu, \sigma^2)}_{p(x_1, \dots, x_n | \gamma, \tau^2, \sigma^2)}$$

$$p(\mu, x_1, \dots, x_n | \gamma, \tau^2, \sigma^2)$$

$$\propto \exp\left\{-\frac{(\mu - \gamma)^2}{2\tau^2}\right\} \prod_n \exp\left\{-\frac{(x_n - \mu)^2}{2\sigma^2}\right\}$$

CONST.

$$= \propto \exp\left\{-\mu^2 \cdot (-) + \mu \cdot (\dots) + \dots\right\} = \mathcal{N}(\mu | \gamma, \tau^2)$$

Posterior of the mean

POSTERIOR: $\mathcal{N}(\mu | \bar{x}, \tau^2)$

$$\bar{x} = \frac{\frac{1}{\tau^2}m + N \cdot \bar{x}}{\frac{1}{\tau^2} + N \cdot \frac{1}{\sigma^2}}$$

$\tau^2 \approx 0$
 $\tau^2 \propto \infty$

$\bar{x}^2 = \frac{1}{1+N} \cdot \frac{1+N}{\frac{1}{\tau^2} + N \cdot \frac{1}{\sigma^2}}$

WEIGHTED HARMONIC AVG of τ^2 & σ^2

POSTERIOR VARIANCE DECREASES $\propto 1/N$ / SAMPLE SIZE

Things to know about Gaussians

- $\hat{\mu}_{MLE} = \frac{1}{N} \sum_n x_n$ $E[\dots] = \mu^2$
 $\hat{\sigma}^2_{MLE} = \frac{1}{N} \sum_n (x_n - \hat{\mu}_{MLE})^2 \parallel \underbrace{\frac{1}{N-1} \sum_n \dots}$
- conjugate for the mean Gaussian
 - variance of posterior decreases with number of samples
 - mean of posterior gets closer and closer to sample average

Gaussians

- modeling singleton continuous variables

NEXT:

- modeling dependence among continuous variables:
 - stock prices given previous few days
 - salary given **GPA**
 - ...

Linear regression *NOT IN x_i ; in w*

GPA SALARY

Input: $(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)$

Goal: learn a map $x \mapsto t$

SALARY

$t = w_0 + w_1 x + w_2 x^2$

LINER in w_0, w_1, w_2

$+ \approx \sum_j w_j \phi_j(x)$

Basis func, FEATURES

$\phi_0(x) = 1$

$\phi_1(x) = \text{GPA}$

$\phi_2(x) = \text{GPA}^2$

$\phi_3(x) = \text{GENDER}$

$\phi_4(x) = \text{GENDER} \cdot \text{GPA}$

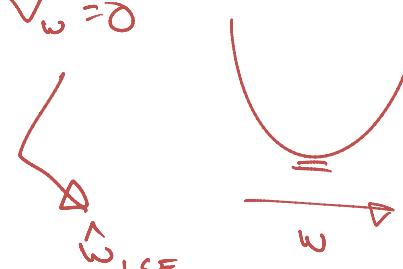
1-PENILES

Linear regression

Input: $(x_1, t_1), \dots, (x_N, t_N)$
 Assume: $t \approx \sum_i w_i \phi_i(x)$

Goal: $\min_w \sum_n (t_n - \underline{w} \cdot \underline{\phi}(x_n))^2$

$\nabla_w = 0$

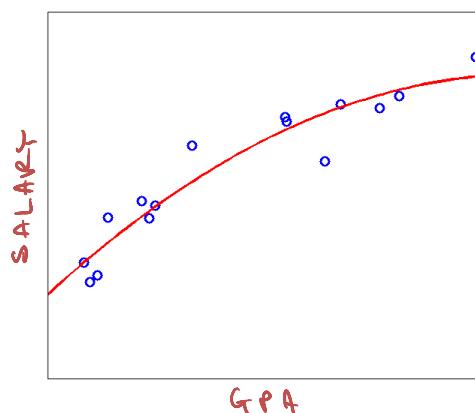


w

Final Prod

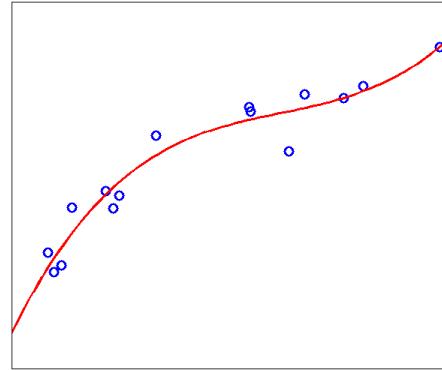
Linear regression

$$t \approx \underline{w}_0 + \underline{w}_1 x + \underline{w}_2 x^2$$



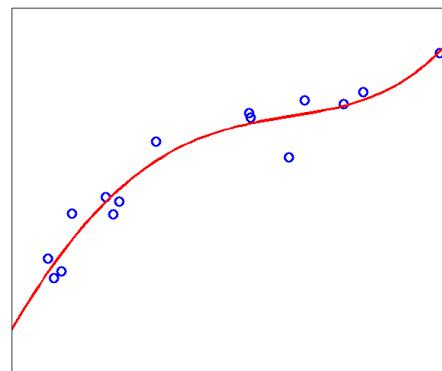
Linear regression

$$t \approx w_0 + w_1 x + w_2 x^2 + w_3 x^3$$



Linear regression

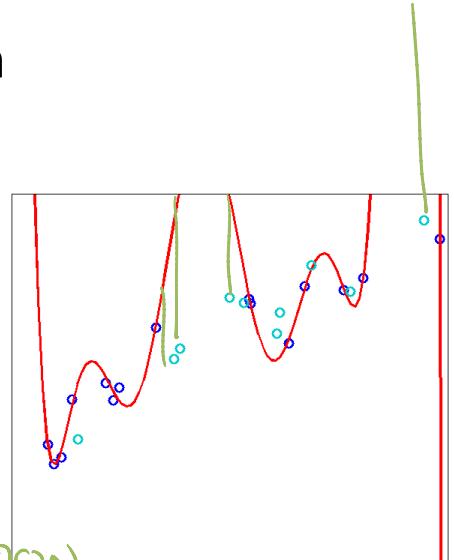
$$t \approx w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$$



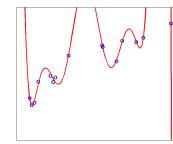
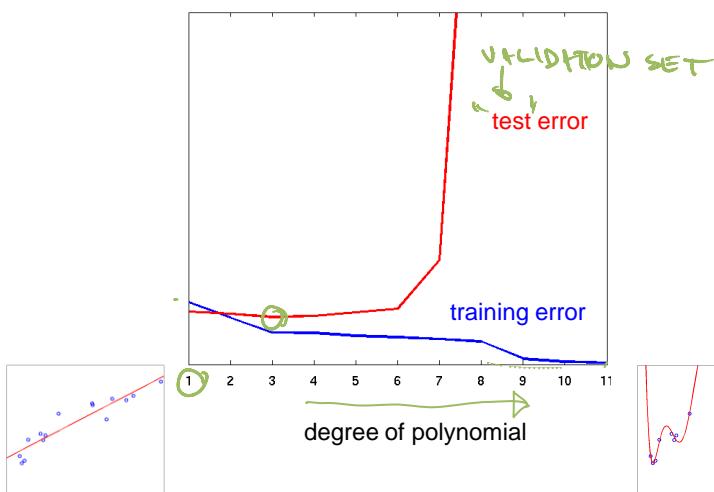
Linear regression

$$t \approx w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + \dots + w_{10} x^{10}$$

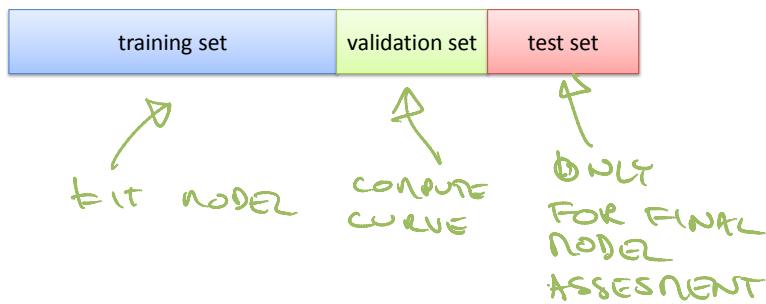
WHAT ABOUT
NEW DATA
 GENERALIZATION



Training error vs Test error



How to select the model?



Project proposals due next WED

- project title
- project idea
- data set to use
- software to write (if any)
- 1-3 papers to read
- teammates
 2-3 people

So far we have covered

Density estimation:

- input: X_1, X_2, \dots, X_N
output: $p(X)$
examples: image denoising/segmentation,
natural language modeling

Regression:

- input: $(x_1, t_1), (x_2, t_2) \dots, (x_N, t_N)$
goal: learn a map $x \mapsto t$
examples: predicting stock value

Soon to cover

Classification:

CATEGORICAL, $\{0,1\}$

- input: $(x_1, y_1), (x_2, y_2) \dots, (x_N, y_N)$
goal: learn a map $x \mapsto y$
examples: object recognition, spam detection

Later in the course

Sequential decisions

- receive observations one at a time
choose action, receive reward
examples: driving a vehicle, controlling robotic arm