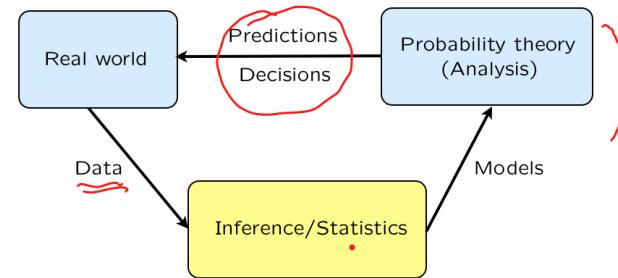


LECTURE 14: Introduction to Bayesian inference

- The big picture
 - motivation, applications
 - problem types (hypothesis testing, estimation, etc.)
- The general framework
 - Bayes' rule → posterior (4 versions)
 - point estimates (MAP, LMS)
 - performance measures (prob. of error; mean squared error)
 - examples

Inference: the big picture



Inference then and now

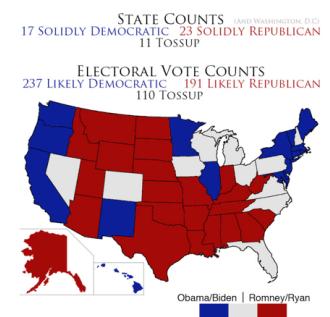
- Then:
10 patients were treated: 3 died
10 patients were not treated: 5 died
Therefore ...

Now:

- Big data
- Big models
- Big computers

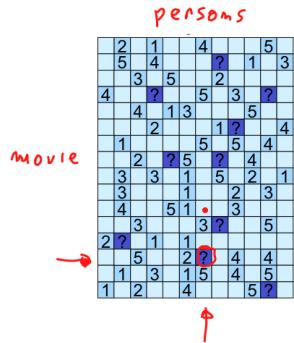
A sample of application domains

- Design and interpretation of experiments
 - polling •



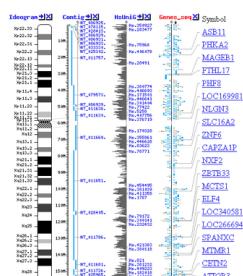
A sample of application domains

- marketing, advertising
- recommendation systems
 - Netflix competition



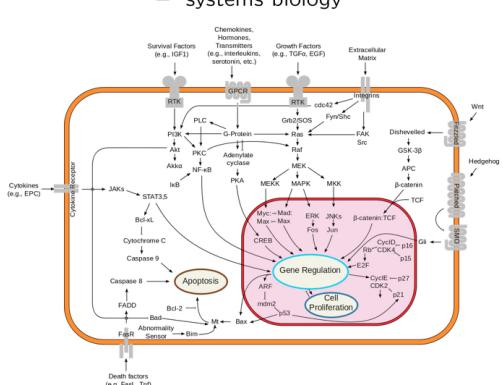
A sample of application domains

- Life sciences
 - genomics



- neuroscience, etc., etc.

– systems biology



A sample of application domains

- Finance



A sample of application domains

- Modeling and monitoring the oceans
- Modeling and monitoring global climate
- Modeling and monitoring pollution
- Interpreting data from physics experiments
- Interpreting astronomy data

A sample of application domains

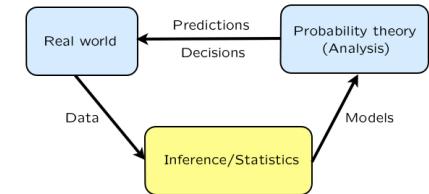
- Signal processing
 - communication systems (noisy ...)
 - speech processing and understanding
 - image processing and understanding
 - tracking of objects
 - positioning systems (e.g., GPS)
 - detection of abnormal events
-

Model building versus inferring unobserved variables



$$X = aS + W$$

- Model building:
 - know "signal" S , observe X
 - infer a
- Variable estimation:
 - know a , observe X
 - infer S



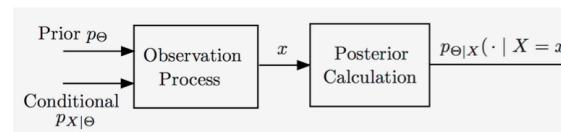
Hypothesis testing versus estimation

- Hypothesis testing:
 - unknown takes one of few possible values
 - aim at small probability of incorrect decision

Is it an airplane or a bird?
- Estimation:
 - numerical unknown(s)
 - aim at an estimate that is "close" to the true but unknown value

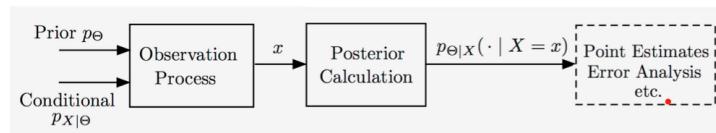
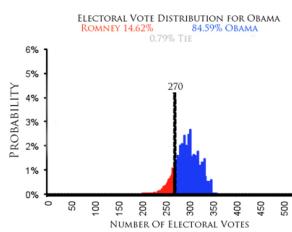
The Bayesian inference framework

- Unknown Θ
 - treated as a random variable
 - prior distribution p_Θ or f_Θ
- Observation X
 - observation model $p_{X|\Theta}$ or $f_{X|\Theta}$
- Use appropriate version of the Bayes rule to find $p_{\Theta|X}(\cdot | X = x)$ or $f_{\Theta|X}(\cdot | X = x)$
- Where does the prior come from?
 - symmetry
 - known range
 - earlier studies
 - subjective or arbitrary



The output of Bayesian inference

The complete answer is a posterior distribution:
PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



Point estimates in Bayesian inference

The complete answer is a posterior distribution:
PMF $p_{\Theta|X}(\cdot | x)$ or PDF $f_{\Theta|X}(\cdot | x)$



estimate: $\hat{\theta} = g(x)$

(number)

estimator: $\widehat{\Theta} = g(X)$

(random variable)

- Maximum a posteriori probability (MAP):

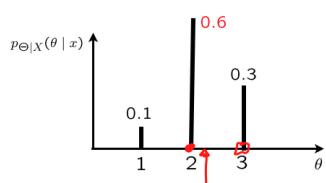
$$p_{\Theta|X}(\theta^* | x) = \max_{\theta} p_{\Theta|X}(\theta | x),$$

$$f_{\Theta|X}(\theta^* | x) = \max_{\theta} f_{\Theta|X}(\theta | x).$$

- Conditional expectation: $E[\Theta | X = x]$ (LMS: Least Mean Squares)

Discrete Θ , discrete X

- values of Θ : alternative hypotheses



$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta) p_{X|\Theta}(x | \theta)}{p_X(x)}$$

$$p_X(x) = \sum_{\theta'} p_\Theta(\theta') p_{X|\Theta}(x | \theta')$$

- conditional prob of error:

$$P(\hat{\theta} \neq \Theta | X = x) = 0.4$$

smallest under the MAP rule

- overall probability of error:

$$\text{MAP rule: } \hat{\theta} = 2$$

LMS: $\hat{\theta} = E[\Theta | X=x] = 2.2$

$$\begin{aligned} P(\hat{\Theta} \neq \Theta) &= \sum_x P(\hat{\Theta} \neq \Theta | X=x) p_X(x) \\ &= \sum_{\theta} P(\hat{\Theta} \neq \Theta | \Theta=\theta) p_\Theta(\theta) \end{aligned}$$

Discrete Θ , continuous X

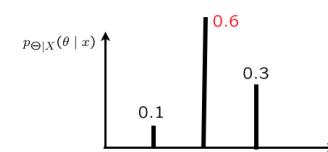
- Standard example:

- send signal $\Theta \in \{1, 2, 3\}$

$$X = \Theta + W$$

$W \sim N(0, \sigma^2)$, indep. of Θ

$$f_{X|\Theta}(x | \theta) = f_W(x - \theta)$$



- MAP rule: $\hat{\theta} = 2$

$$p_{\Theta|X}(\theta | x) = \frac{p_\Theta(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

$$f_X(x) = \sum_{\theta'} p_\Theta(\theta') f_{X|\Theta}(x | \theta')$$

- conditional prob of error:

$$P(\hat{\theta} \neq \Theta | X = x)$$

smallest under the MAP rule

- overall probability of error:

$$\begin{aligned} P(\hat{\Theta} \neq \Theta) &= \int P(\hat{\Theta} \neq \Theta | X=x) f_X(x) dx \\ &= \sum_{\theta} P(\hat{\Theta} \neq \Theta | \Theta=\theta) p_\Theta(\theta) \end{aligned}$$

Continuous Θ , continuous X

- linear normal models
estimation of a noisy signal
 $X = \Theta + W$

Θ and W : independent normals

multi-dimensional versions (many normal parameters, many observations)

- estimating the parameter of a uniform
 X : uniform $[0, \Theta]$
 Θ : uniform $[0, 1]$

$$\hat{\Theta} = g(X) \quad \begin{matrix} MAP \\ LMS \end{matrix}$$

interested in:

$$\left\{ \begin{array}{l} E[(\hat{\Theta} - \Theta)^2 | X = x] \\ E[(\hat{\Theta} - \Theta)^2] \end{array} \right.$$

$$f_{\Theta|X}(\theta | x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

$$f_X(x) = \int f_{\Theta}(\theta') f_{X|\Theta}(x | \theta') d\theta'$$

Inferring the unknown bias of a coin and the Beta distribution

- Standard example:
 - coin with bias Θ ; prior $f_{\Theta}(\cdot)$
 - fix n ; K = number of heads
- Assume $f_{\Theta}(\cdot)$ is uniform in $[0, 1]$

$$f_{\Theta|K}(\theta | k) = \frac{1}{P_K(k)} \theta^k (1-\theta)^{n-k} \quad \theta \in [0, 1]$$

$$= \frac{1}{d(n, k)} \theta^k (1-\theta)^{n-k} \quad \text{"Beta distribution, with parameters } (k+1, n-k+1)\text{"}$$

- If prior is Beta: $f_{\Theta}(\theta) = \frac{1}{c} \theta^\alpha (1-\theta)^\beta \quad \alpha, \beta \geq 0$

$$f_{\Theta|K}(\theta | k) = \frac{1}{c} \theta^\alpha (1-\theta)^\beta \binom{n}{k} \theta^k (1-\theta)^{n-k} / P_K(k) = d \theta^{\alpha+k} (1-\theta)^{\beta+n-k}$$

Inferring the unknown bias of a coin: point estimates

- Standard example:
 - coin with bias Θ ; prior $f_{\Theta}(\cdot)$
 - fix n ; K = number of heads

- Assume $f_{\Theta}(\cdot)$ is uniform in $[0, 1]$

$$f_{\Theta|K}(\theta | k) = \frac{1}{d(n, k)} \theta^k (1-\theta)^{n-k}$$

- MAP estimate:

$$\hat{\theta}_{MAP} = \frac{k}{n}$$

$$\max_{\theta} [k \log \theta + (n-k) \log (1-\theta)]$$

$$k/\theta \rightarrow (n-k)/(1-\theta) = 0$$

$$\hat{\theta}_{MAP} = \frac{k}{n}$$

$$E[\Theta | K = k] = \int_0^1 \theta f_{\Theta|K}(\theta | k) d\theta \quad \alpha \geq 0, \beta \geq 0$$

$$= \frac{1}{d(n, k)} \int_0^1 \theta^{k+1} (1-\theta)^{n-k} d\theta$$

$$= \frac{1}{\binom{n}{k}} \cdot \frac{(k+1)! (n-k)!}{(n+2)!}$$

$$= \frac{k+1}{n+2} \approx \frac{k}{n} \quad (\text{large } n)$$

Summary

- Problem data: $p_{\Theta}(\cdot)$, $p_{X|\Theta}(\cdot | \cdot)$
- Given the value x of X : find, e.g., $p_{\Theta|X}(\cdot | x)$
 - using appropriate version of the Bayes rule (4 choices)
- Estimator $\hat{\Theta} = g(X)$ Estimate $\hat{\theta} = g(x)$
 - MAP: $\hat{\theta}_{MAP} = g_{MAP}(x)$ maximizes $p_{\Theta|X}(\theta | x)$
 - LMS: $\hat{\theta}_{LMS} = g_{LMS}(x) = E[\Theta | X = x]$
- Performance evaluation of an estimator $\hat{\Theta}$
 - $P(\hat{\Theta} \neq \Theta | X = x)$ $\rightarrow E[(\hat{\Theta} - \Theta)^2 | X = x]$
 - $P(\hat{\Theta} \neq \Theta)$ $\rightarrow E[(\hat{\Theta} - \Theta)^2] \cdot \text{total prob exp}$

LECTURE 15: Linear models with normal noise

$$X_i = \sum_{j=1}^m a_{ij} \Theta_j + W_i \quad W_i, \Theta_j: \text{independent, normal}$$

- Very common and convenient model
- Bayes' rule: normal posteriors
- MAP and LMS estimates coincide
 - simple formulas
(linear in the observations)
- Many nice properties
- Trajectory estimation example

Recognizing normal PDFs

$$X \sim N(\mu, \sigma^2) \quad f_X(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2} \quad 2\alpha x + \beta = 0$$

$$c \cdot e^{-8(x-3)^2} \quad \mu = 3 \quad \frac{1}{2\sigma^2} = 8 \Rightarrow \sigma^2 = \frac{1}{16} \quad c = \frac{1}{\frac{1}{4}\sqrt{2\pi}}$$

$$f_X(x) = c \cdot e^{-(\alpha x^2 + \beta x + \gamma)} \quad \alpha > 0 \quad \text{Normal with mean } -\beta/2\alpha \text{ and variance } 1/2\alpha$$

$$\alpha x^2 + \beta x + \gamma = \alpha \left(x^2 + \frac{\beta}{\alpha} x + \frac{\gamma}{\alpha} \right) = \alpha \left(\left(x + \frac{\beta}{2\alpha} \right)^2 - \frac{\beta^2}{4\alpha^2} + \frac{\gamma}{\alpha} \right)$$

$$f_X(x) = c \underbrace{e^{-\alpha \left(x + \frac{\beta}{2\alpha} \right)^2}}_{\text{Normal}} \underbrace{e^{-\alpha \left(-\frac{\beta^2}{4\alpha^2} + \frac{\gamma}{\alpha} \right)}}_{\text{constant}} \quad \mu = -\frac{\beta}{2\alpha}$$

$$\frac{1}{2\sigma^2} = \alpha \Rightarrow \sigma^2 = 1/2\alpha$$

Estimating a normal random variable in the presence of additive normal noise

$$X = \Theta + W \quad \Theta, W : N(0, 1), \text{ independent}$$

$$f_{X|\Theta}(x|\theta) : X = \theta + W \quad N(\theta, 1)$$

$$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x|\theta)}{f_X(x)} = \frac{c e^{-\frac{1}{2}\theta^2} c e^{-\frac{1}{2}(x-\theta)^2}}{f_X(x)} = \underline{c(x)} e^{-\text{quadratic}(\theta)}$$

$$\text{Fix } x \quad \min_{\theta} \left[\frac{1}{2}\theta^2 + \frac{1}{2}(x-\theta)^2 \right] \quad \theta + (x-\theta) = 0$$

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = \mathbb{E}[\Theta | X = x] = \underline{x/2}$$

$$\hat{\Theta}_{\text{MAP}} = \mathbb{E}[\Theta | X] = \underline{x/2}$$

Estimating a normal parameter in the presence of additive normal noise

$$X = \Theta + W \quad \Theta, W : N(0, 1), \text{ independent}$$

$$f_{\Theta|X}(\theta|x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x|\theta)}{f_X(x)}$$

$$f_X(x) = \int f_{\Theta}(\theta) f_{X|\Theta}(x|\theta) d\theta$$

$$\hat{\Theta}_{\text{MAP}} = \hat{\Theta}_{\text{LMS}} = \mathbb{E}[\Theta | X] = \frac{X}{2}$$

- Even with general means and variances:

- posterior is normal
- LMS and MAP estimators coincide
- these estimators are “linear,” of the form $\hat{\Theta} = aX + b$

The case of multiple observations

$$\begin{aligned} X_1 &= \Theta + W_1 & \Theta &\sim N(x_0, \sigma_0^2) & W_i &\sim N(0, \sigma_i^2) \\ &\vdots \\ X_n &= \Theta + W_n & \Theta, W_1, \dots, W_n &\text{ independent} \end{aligned}$$

$$f_{\Theta|X}(\theta | x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

$$f_X(x) = \int f_{\Theta}(\theta) f_{X|\Theta}(x | \theta) d\theta$$

$$f_{X_i|\Theta}(x_i | \theta) = c_i e^{-(x_i - \theta)^2 / 2\sigma_i^2}$$

given $\Theta = \theta$: $X_i = \theta + W_i \sim N(\theta, \sigma_i^2)$

$$f_{X|\Theta}(x | \theta) = f_{x_1, \dots, x_n | \Theta}(x_1, \dots, x_n | \theta) = \prod_{i=1}^n f_{x_i | \Theta}(x_i | \theta)$$

given $\Theta = \theta$: W_i independent $\Rightarrow X_i$ independent

$$f_{\Theta|X}(\theta | x) = \frac{1}{f_X(x)} \cdot c_0 e^{-(\theta - x_0)^2 / 2\sigma_0^2} \prod_{i=1}^n c_i e^{-(x_i - \theta)^2 / 2\sigma_i^2} \quad \text{Normal!}$$

The case of multiple observations

- Key conclusions:
 - posterior is normal
 - LMS and MAP estimates coincide
 - these estimates are "linear," of the form $\hat{\theta} = a_0 + a_1 x_1 + \dots + a_n x_n$
- Interpretations:
 - estimate $\hat{\theta}$: weighted average of x_0 (prior mean) and x_i (observations)
 - weights determined by variances

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = \mathbb{E}[\Theta | X = x] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2} \cdot \frac{1}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

σ_i^2 large

x_i very noisy
 \Rightarrow small weight

The case of multiple observations

$$f_{\Theta|X}(\theta | x) = c \cdot \exp \{-\text{quad}(\theta)\} \quad \text{quad}(\theta) = \frac{(\theta - x_0)^2}{2\sigma_0^2} + \frac{(\theta - x_1)^2}{2\sigma_1^2} + \dots + \frac{(\theta - x_n)^2}{2\sigma_n^2}$$

find peak

$$\frac{d}{d\theta} \text{quad}(\theta) = 0: \sum_{i=0}^n \frac{(\theta - x_i)}{\sigma_i^2} = 0 \Rightarrow \theta \sum_{i=0}^n \frac{1}{\sigma_i^2} = \sum_{i=0}^n \frac{x_i}{\sigma_i^2}$$

$$\hat{\theta}_{\text{MAP}} = \hat{\theta}_{\text{LMS}} = \mathbb{E}[\Theta | X = x] = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

The mean squared error

$$f_{\Theta|X}(\theta | x) = c \cdot \exp \{-\text{quad}(\theta)\}$$

$$\text{quad}(\theta) = \frac{(\theta - x_0)^2}{2\sigma_0^2} + \frac{(\theta - x_1)^2}{2\sigma_1^2} + \dots + \frac{(\theta - x_n)^2}{2\sigma_n^2}$$

$$\hat{\theta} = \frac{\sum_{i=0}^n \frac{x_i}{\sigma_i^2}}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

Performance measures:

$$\mathbb{E}[(\Theta - \hat{\Theta})^2 | X = x] = \mathbb{E}[(\Theta - \hat{\theta})^2 | X = x] = \text{var}(\Theta | X = x) = \frac{1}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

$$\mathbb{E}[(\Theta - \hat{\Theta})^2] = \int \mathbb{E}[(\Theta - \hat{\theta})^2 | X = x] f_X(x) dx$$

$$f_X(x) = c \cdot e^{-(\alpha x^2 + \beta x + \gamma)} \quad \alpha > 0 \quad \text{Normal with mean } -\beta/2\alpha \text{ and variance } 1/2\alpha$$

$$\alpha = \frac{1}{2\sigma_0^2} + \dots + \frac{1}{2\sigma_n^2}$$

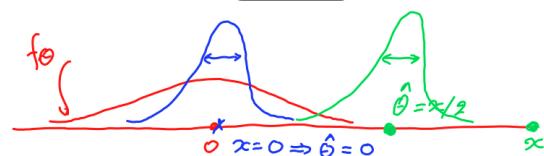
Some σ_i^2 small \rightarrow MSE small
all σ_i^2 large \rightarrow MSE large

The mean squared error

$$\mathbb{E}[(\Theta - \hat{\Theta})^2 | X = x] \stackrel{?}{=} \mathbb{E}[(\Theta - \hat{\Theta})^2] = 1 / \sum_{i=0}^n \frac{1}{\sigma_i^2}$$

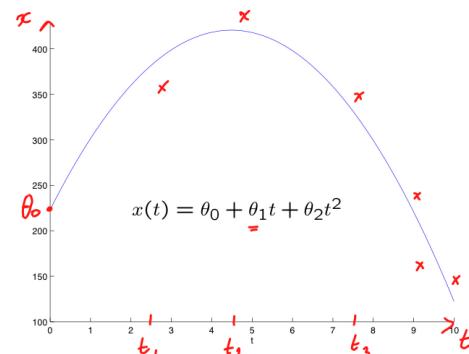
- Example: $\sigma_0^2 = \sigma_1^2 = \dots = \sigma_n^2 = \sigma^2$ $\frac{1}{(n+1)\sigma^2} = \frac{\sigma^2}{n+1}$
- conditional mean squared error same for all x

- Example: $X = \Theta + W$ $\Theta \sim N(0, 1)$, $W \sim N(0, 1)$
independent Θ, W $\hat{\Theta} = X/2$ $\mathbb{E}[(\Theta - \hat{\Theta})^2 | X = x] = 1/2$



$$\hat{\theta} = \frac{\sum_{i=0}^n x_i}{\sum_{i=0}^n \frac{1}{\sigma_i^2}}$$

The case of multiple parameters: trajectory estimation



- Random variables $\Theta_0, \Theta_1, \Theta_2$ independent; priors f_{Θ_j}
- Measurements at times t_1, \dots, t_n
 $X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$
noise model: f_{W_i}
independent W_i ; independent from Θ_j

A model with normality assumptions

$$X_i = \underline{\Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i} \quad i = 1, \dots, n$$

$$f_{\Theta|X}(\theta | x) = \frac{f_{\Theta}(\theta) f_{X|\Theta}(x | \theta)}{f_X(x)}$$

$$f_X(x) = \int f_{\Theta}(\theta) f_{X|\Theta}(x | \theta) d\theta$$

- assume $\Theta_j \sim N(0, \sigma_j^2)$, $W_i \sim N(0, \sigma^2)$; independent
- Given $\Theta = \theta = (\theta_0, \theta_1, \theta_2)$, X_i is: $N(\theta_0 + \theta_1 t_i + \theta_2 t_i^2, \sigma^2)$

$$f_{X_i|\Theta}(x_i | \theta) = c \cdot \exp \left\{ -\frac{(x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2}{2\sigma^2} \right\}$$

- posterior: $f_{\Theta|X}(\theta | x) = \frac{1}{f_X(x)} \prod_{j=0}^2 f_{\Theta_j}(\theta_j) \prod_{i=1}^n f_{X_i|\Theta}(x_i | \theta)$

$$c(x) \exp \left\{ -\frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 \right\}$$

A model with normality assumptions

$$f_{\Theta|X}(\theta | x) = \underline{c(x) \exp \left\{ -\frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 \right\}}$$

- MAP estimate: maximize over $(\theta_0, \theta_1, \theta_2)$;
(minimize quadratic function)

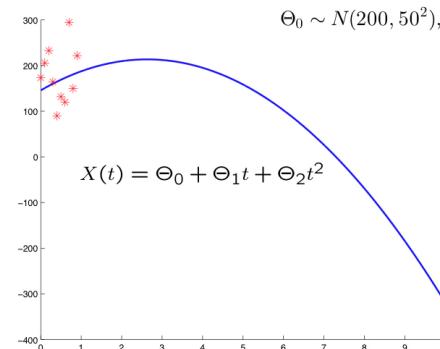
$$\frac{\partial}{\partial \theta_j} (\text{quad}(\theta)) = 0 \quad \begin{matrix} 3 \text{ equations, } 3 \text{ unknowns} \\ \uparrow \text{linear} \end{matrix}$$

Linear normal models .

- Θ_j and X_i are linear functions of independent normal random variables
- $f_{\Theta|X}(\theta|x) = c(x) \exp\{-\text{quadratic}(\theta_1, \dots, \theta_m)\}$
- MAP estimate: maximize over $(\theta_1, \dots, \theta_m)$; linear regression
- $\widehat{\Theta}_{MAP,j}$: linear function of $X = (X_1, \dots, X_n)$
- Facts:
 - $\widehat{\Theta}_{MAP,j} = E[\Theta_j | X]$
 - marginal posterior PDF of Θ_j : $f_{\Theta_j|X}(\theta_j | x)$, is normal
 - MAP estimate based on the joint posterior PDF:
same as MAP estimate based on the marginal posterior PDF
 - $E[(\widehat{\Theta}_{i,MAP} - \Theta_i)^2 | X = x]$: same for all x

An illustration

Estimating the trajectory of a free-falling object



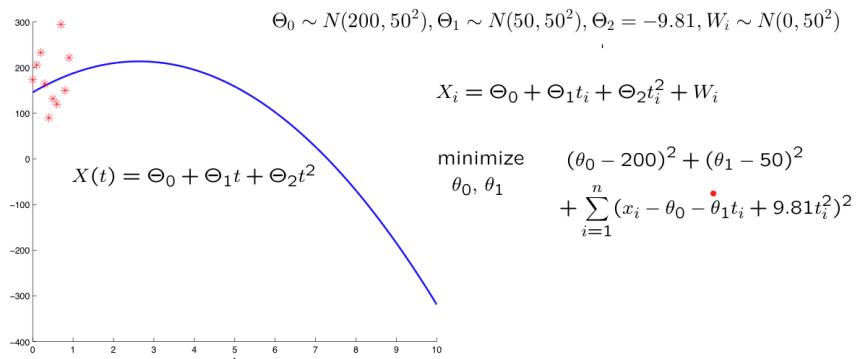
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \underset{\theta_0, \theta_1, \theta_2}{\text{minimize}} \quad \frac{1}{2} \left(\frac{\theta_0^2}{\sigma_0^2} + \frac{\theta_1^2}{\sigma_1^2} + \frac{\theta_2^2}{\sigma_2^2} \right) \\ & + \frac{1}{\sigma^2} \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i - \theta_2 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



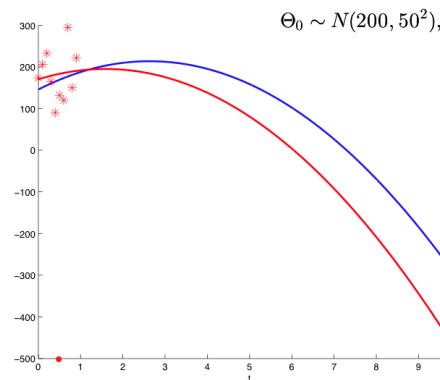
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X(t) = \Theta_0 + \Theta_1 t + \Theta_2 t^2$$

$$\begin{aligned} & \underset{\theta_0, \theta_1}{\text{minimize}} \quad (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

An illustration

Estimating the trajectory of a free-falling object



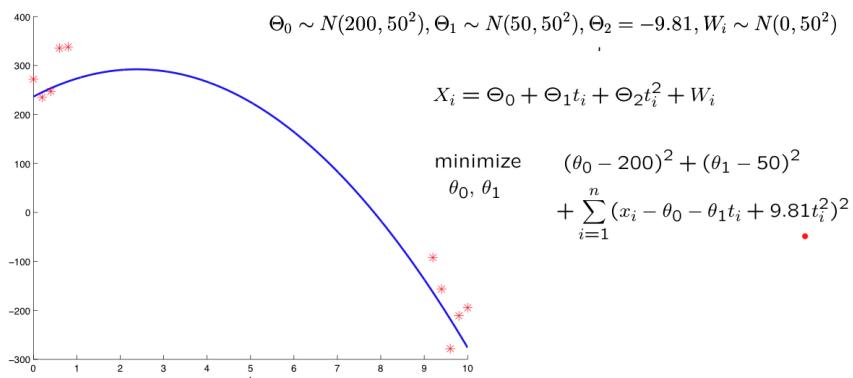
$$\Theta_0 \sim N(200, 50^2), \Theta_1 \sim N(50, 50^2), \Theta_2 = -9.81, W_i \sim N(0, 50^2)$$

$$X_i = \Theta_0 + \Theta_1 t_i + \Theta_2 t_i^2 + W_i$$

$$\begin{aligned} & \underset{\theta_0, \theta_1}{\text{minimize}} \quad (\theta_0 - 200)^2 + (\theta_1 - 50)^2 \\ & + \sum_{i=1}^n (x_i - \theta_0 - \theta_1 t_i + 9.81 t_i^2)^2 \end{aligned}$$

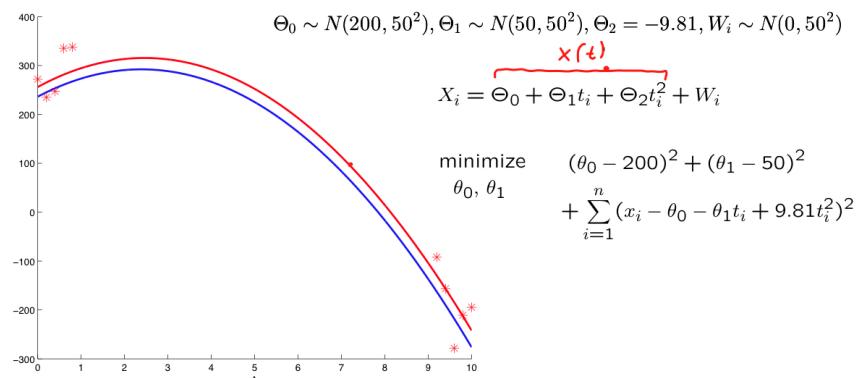
An illustration

Estimating the trajectory of a free-falling object



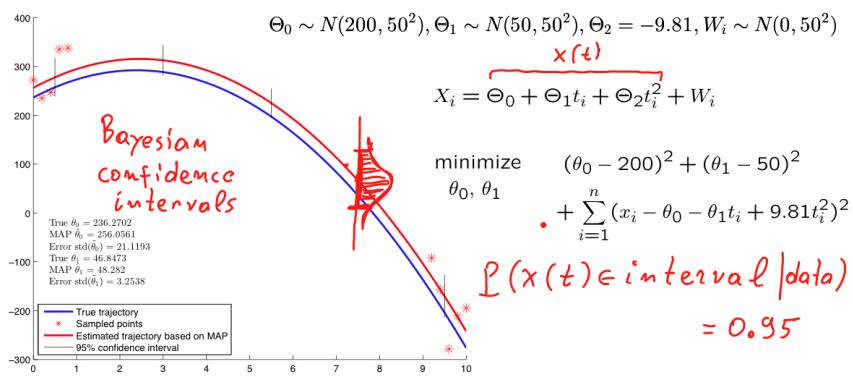
An illustration

Estimating the trajectory of a free-falling object



An illustration

Estimating the trajectory of a free-falling object

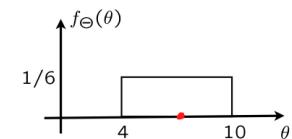


LECTURE 16: Least mean squares (LMS) estimation

- minimize (conditional) mean squared error $E[(\Theta - \hat{\theta})^2 | X = x]$
 - solution: $\hat{\theta} = E[\Theta | X = x]$
 - general estimation method
- Mathematical properties
- Example

LMS estimation in the absence of observations

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
 - no observations available
 - MAP rule: $\text{arg } \hat{\theta} \in [4, 10]$
 - (Conditional) expectation: $\hat{\theta} = 7$
- Criterion: Mean Squared Error (MSE): $E[(\Theta - \hat{\theta})^2]$
minimize mean squared error



LMS estimation in the absence of observations

- Least mean squares formulation:

$$\begin{aligned} &\text{minimize mean squared error (MSE), } E[(\Theta - \hat{\theta})^2]: \hat{\theta} = E[\Theta] . \\ &E[\Theta^2] - 2E[\Theta]\hat{\theta} + \hat{\theta}^2 \quad \frac{d}{d\hat{\theta}} = 0 : -2E[\Theta] + 2\hat{\theta} = 0 \\ &\hat{\theta} = E[\Theta] \end{aligned}$$

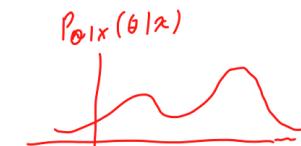
$$\text{Var}(\Theta - \hat{\theta}) + (E[\Theta - \hat{\theta}])^2$$

minimized
when $\hat{\theta} = E[\Theta]$

- Optimal mean squared error: $E[(\Theta - E[\Theta])^2] = \text{var}(\Theta)$

LMS estimation of Θ based on X

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
- observation X ; model $p_{X|\Theta}(x|\theta)$
 - observe that $X = x$



$$\text{minimize mean squared error (MSE), } E[(\Theta - \hat{\theta})^2]: \hat{\theta} = E[\Theta]$$

$$\text{minimize conditional mean squared error, } E[(\Theta - \hat{\theta})^2 | X = x]: \hat{\theta} = E[\Theta | X = x]$$

- LMS estimate: $\hat{\theta} = E[\Theta | X = x]$

estimator: $\hat{\Theta} = E[\Theta | X]$

LMS estimation of Θ based on X

- $E[\Theta]$ minimizes $E[(\Theta - \hat{\theta})^2]$

$$E[(\Theta - E[\Theta])^2] \leq E[(\Theta - c)^2], \text{ for all } c$$

- $E[\Theta | X = x]$ minimizes $E[(\Theta - \hat{\theta})^2 | X = x]$

$$E[(\Theta - E[\Theta | X = x])^2 | X = x] \leq E[(\Theta - g(x))^2 | X = x] \text{ for all } x$$

$$E[(\Theta - E[\Theta | X])^2 | X] \leq E[(\Theta - g(x))^2 | X]$$

$$E[(\Theta - E[\Theta | X])^2] \leq E[(\Theta - g(x))^2]$$



$\hat{\Theta}_{LMS} = E[\Theta | X]$ minimizes $E[(\Theta - g(X))^2]$, over all estimators $\hat{\Theta} = g(X)$

LMS performance evaluation

- LMS estimate: $\hat{\theta} = E[\Theta | X = x]$

estimator: $\hat{\Theta} = E[\Theta | X]$

- Expected performance, once we have a measurement:

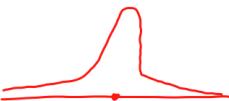
$$\text{MSE} = E[(\Theta - E[\Theta | X = x])^2 | X = x] = \underline{\text{var}(\Theta | X = x)}$$

- Expected performance of the design:

$$\text{MSE} = E[(\Theta - E[\Theta | X])^2] = E[\underline{\text{var}(\Theta | X)}]$$

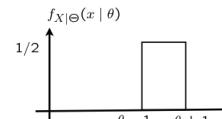
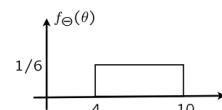
LMS estimation of Θ based on X

- LMS relevant to estimation (not hypothesis testing)

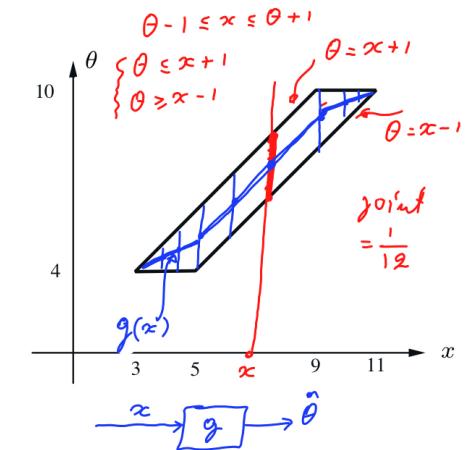


- Same as MAP if the posterior is unimodal and symmetric around the mean
 - e.g., when posterior is normal (the case in "linear-normal" models)

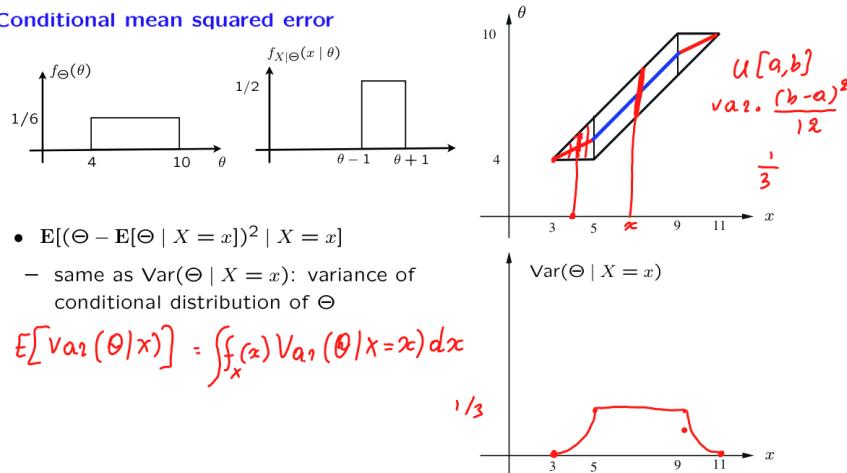
Example



$$x = \theta + u \quad u \sim \text{unif}(-1, 1)$$



Conditional mean squared error



- $E[(\Theta - E[\Theta | X=x])^2 | X=x]$
 - same as $\text{Var}(\Theta | X=x)$: variance of conditional distribution of Θ

$$E[\text{Var}(\Theta | X=x)] = \int_x f(x) \text{Var}(\Theta | X=x) dx$$

LMS estimation with multiple observations or unknowns

- unknown Θ ; prior $p_\Theta(\theta)$
 - interested in a point estimate $\hat{\theta}$
- observations $X = (X_1, X_2, \dots, X_n)$; model $p_{X|\Theta}(x|\theta)$
 - observe that $X = x$
 - new universe: condition on $X = x$
- LMS estimate: $E[\Theta | X_1 = x_1, \dots, X_n = x_n]$
- If Θ is a vector, apply to each component separately

$$\Theta = (\theta_1, \dots, \theta_m) \quad \hat{\Theta}_j = E[\Theta_j | X_1 = x_1, \dots, X_n = x_n]$$

Some challenges in LMS estimation

- Full correct model, $f_{X|\Theta}(x|\theta)$, may not be available
- Can be hard to compute/implement/analyze

$$E[\Theta_j | X=x] = \iiint \theta_j f_{\Theta|X}(\theta|x) d\theta_1 \dots d\theta_m$$

$$f_{\Theta|X}(\theta|x) = \frac{f_\Theta(\theta) f_{X|\Theta}(x|\theta)}{f_X(x)}$$

$$f_X(x) = \int f_\Theta(\theta') f_{X|\Theta}(x|\theta') d\theta'$$

Properties of the estimation error in LMS estimation

- Estimator: $\hat{\Theta} = E[\Theta | X]$
- Error: $\bar{\Theta} = \hat{\Theta} - \Theta$

$$E[\bar{\Theta} | X=x] = 0$$

$$E[\hat{\Theta} - \Theta | X=x] = \hat{\Theta} - E[\Theta | X=x] = 0$$

$$\text{cov}(\bar{\Theta}, \bar{\Theta}) = 0$$

$$E[\bar{\Theta} \hat{\Theta}] - E[\bar{\Theta}] E[\hat{\Theta}] = 0$$

$$E[\bar{\Theta} \hat{\Theta} | X=x] = \hat{\Theta} E[\bar{\Theta} | X=x] = 0$$

$$\text{var}(\Theta) = \text{var}(\hat{\Theta}) + \text{var}(\bar{\Theta})$$

$$\Theta = \hat{\Theta} - \bar{\Theta}$$