

**6.s077 — INTRODUCTION TO DATA SCIENCE
EECS, MIT, Spring 2018**

Lecture 5

**Multidimensional data, $\mathbb{E}[Y | X]$
Bayesian Inference: posterior distribution**

Today's agenda — Part I

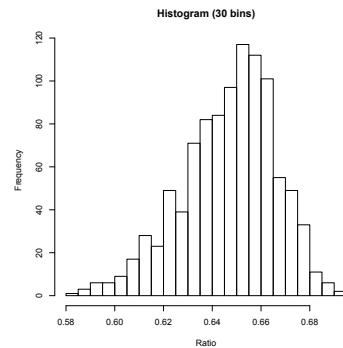
- Multidimensional data
 - visualization
 - i.i.d. models
 - estimation (plugin, feature matching, ML)
- Predicting one variable from another
 - model-based, $\mathbb{E}[Y | X]$: today
 - based on data (X_i, Y_i) : come back during regression lectures

Multidimensional data

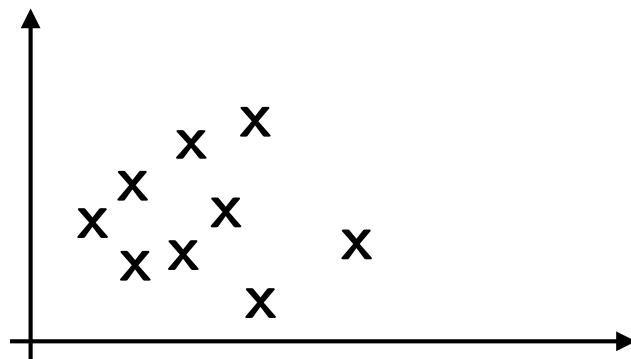
- i th sample (x_i, y_i, z_i)
e.g., (gender, height, weight) of i th person

- Visualize:

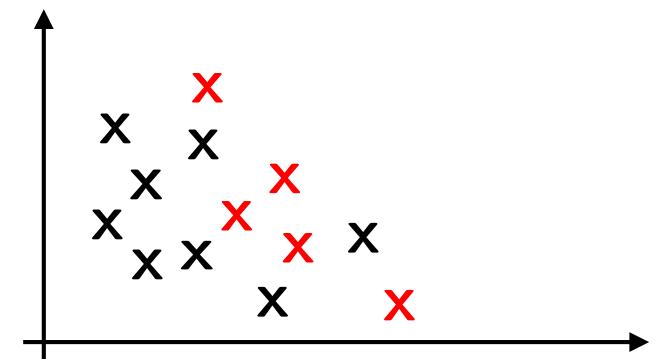
one at a time



two at a time



three, one categorical



“Standard model:” Independent random sampling

- Assume two-dimensional data (only for exposition purposes)
- \mathbb{P} : distribution on \mathbb{R}^2 (discrete; continuous; mixed)

$$p_{X,Y}(x,y) \quad f_{X,Y}(x,y)$$

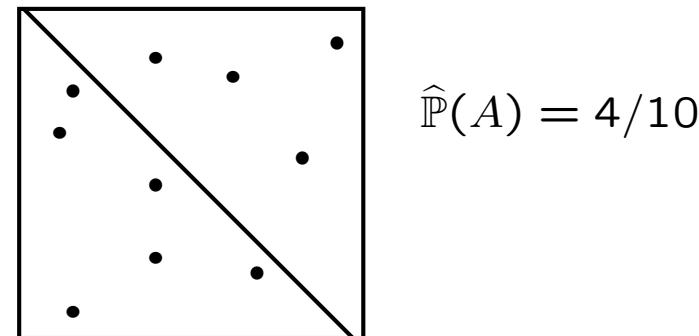
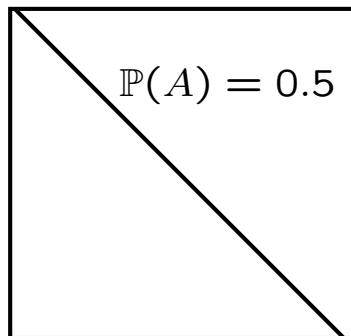
- $(X_i, Y_i) \sim \mathbb{P}$, independent for different i

X_i and Y_i usually dependent

- $\hat{\mathbb{P}}$: empirical distribution on \mathbb{R}^2
 - weight $1/n$ on each data point

$\hat{\mathbb{P}}(A) =$ fraction of data that belong to A

uniform \mathbb{P}



Parameter estimation with multidimensional data

- Nothing new!

- **Plugin:**

$$\theta = \mathbb{E}[g(X, Y)] \quad \widehat{\Theta} = \widehat{\mathbb{E}}[g(X, Y)] = \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i)$$

$$\theta = \text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

$$\widehat{\Theta} = \widehat{\mathbb{E}}[XY] - \widehat{\mathbb{E}}[X]\widehat{\mathbb{E}}[Y]$$

Parameter estimation with multidimensional data (ctd.)

- Feature matching and the method of moments
- Example:

$$\mathbb{P}^\theta \quad \theta = (\theta_1, \theta_2, \theta_3)$$

Choose $\mathbb{E}^\theta[X]$, $\mathbb{E}^\theta[Y]$, $\mathbb{E}^\theta[XY]$ as features

$$\mathbb{E}^\theta[X] = \hat{\mathbb{E}}[X]$$

$$\mathbb{E}^\theta[Y] = \hat{\mathbb{E}}[Y] \qquad \text{solve for } \theta = (\theta_1, \theta_2, \theta_3)$$

$$\mathbb{E}^\theta[XY] = \hat{\mathbb{E}}[XY]$$

- easiest when we have closed-form expressions
(as a function of θ) for left-hand side

Parameter estimation with multidimensional data (ctd.)

- Maximum Likelihood

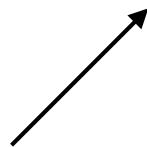
$$\max_{\theta} \prod_{i=1}^n \mathbb{P}^{\theta}(X_i = x_i) \quad \prod_{i=1}^n f_X^{\theta}(x_i)$$

X_i, x_i are now vectors

f_X is a joint PDF

- Nothing new
 - nice theoretical properties remain

- Sampling distribution, bootstrap, confidence intervals:
again nothing new



one parameter at a time

What is new? New types of prediction problems

$$\begin{array}{c} \theta? \quad \mathbb{P}_Y^\theta \\ \hline Y_1 \\ \vdots \\ Y_n \\ \hline Y? \end{array}$$

i.i.d.

$$\min_{\hat{y}} \quad \mathbb{E}^\theta \left[(Y - \hat{y})^2 \right]$$

$$\begin{array}{c} \theta? \quad \mathbb{P}_{X,Y}^\theta \\ \hline X_1 \quad | \quad Y_1 \\ \vdots \quad | \quad \vdots \\ X_n \quad | \quad Y_n \\ \hline X \quad | \quad Y? \end{array}$$

$$\min_{g(\cdot)} \quad \mathbb{E}^\theta \left[(Y - g(X))^2 \right]$$

$\hat{Y} = g(X)$ estimator

A. Have \mathbb{P}_Y^θ , know θ

$$\hat{Y} = \mathbb{E}^\theta[Y]$$

Have $\mathbb{P}_{X,Y}^\theta$, know θ

$$\hat{Y} = \mathbb{E}^\theta[Y | X]$$

data $X_1, Y_1, \dots, X_n, Y_n$ are irrelevant

Some details

$$\mathbb{E}[(Y - \mathbb{E}[Y])^2] \leq \mathbb{E}[(Y - c)^2] \quad \text{for all } c$$

apply to conditional universe, where $X = x$ has occurred

$$\mathbb{E}[(Y - \mathbb{E}[Y | X = x])^2 | X = x] \leq \mathbb{E}[(Y - c)^2 | X = x], \quad \text{for all } c$$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X = x])^2 | X = x] \leq \mathbb{E}[(Y - g(x))^2 | X = x], \quad \text{for all } g(\cdot)$$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X])^2 | X] \leq \mathbb{E}[(Y - g(X))^2 | X], \quad \text{for all } g(\cdot)$$

take expectation of both sides; use law of iterated expectations,
and the property $X \leq Y \Rightarrow \mathbb{E}[X] \leq \mathbb{E}[Y]$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X])^2] \leq \mathbb{E}[(Y - g(X))^2], \quad \text{for all } g(\cdot)$$

Prediction via a model versus directly from data

B. Have $\mathbb{P}_{X,Y}^\theta$, do not know θ

$\theta?$	$\mathbb{P}_{X,Y}^\theta$
X_1	Y_1
:	:
:	:
X_n	Y_n
<hr/>	
X	$Y?$

→ regression

Today's agenda — Part II

- The Bayesian framework → posterior distributions
- Four versions of Bayes' rule
- Bayesian confidence intervals
- Coin parameter example

Today's agenda — Part II

$\theta?$ \mathbb{P}_X^θ

X_1

\vdots

\vdots

X_n

- Classical/frequentist statistics: θ : unknown constant
- Bayesian philosophy
 - parameter treated as (unobserved) realized value of **random variable** Θ
 - prior distribution \mathbb{P}_Θ : $p_\Theta(x)$ or $f_\Theta(\theta)$
joint PMF/PDF if θ is a vector
- Observation model $\mathbb{P}_{X|\Theta}$: $p_{X|\Theta}(x|\theta)$ or $f_{X|\Theta}(x|\theta)$
replace $p^\theta(x)$, $f^\theta(x)$
- Prior comes from: symmetry; known range;
earlier studies; subjective beliefs

The big picture

$$\frac{\theta? \quad \mathbb{P}_X^\theta}{X_1 \quad \vdots \quad X_n}$$

Classical
inference

$$\frac{\mathbb{P}_X \quad \mathbb{P}_{Y|X}}{X \quad | \quad Y?}$$

Model-based
prediction

$$\mathbb{E}[Y | X]$$

$$\frac{\mathbb{P}_\Theta \quad \mathbb{P}_{X|\Theta}}{\Theta? \quad | \quad X}$$

Bayesian
inference

$$\mathbb{E}[\Theta | X]$$

$$\mathbb{P}_{\Theta|X} = \frac{\mathbb{P}_\Theta \cdot \mathbb{P}_{X|\Theta}}{\mathbb{P}_X}$$

Two different flavors

$$\frac{\mathbb{P}_{U,V,W,Y}}{(U, V, W) \quad | \quad Y?}$$

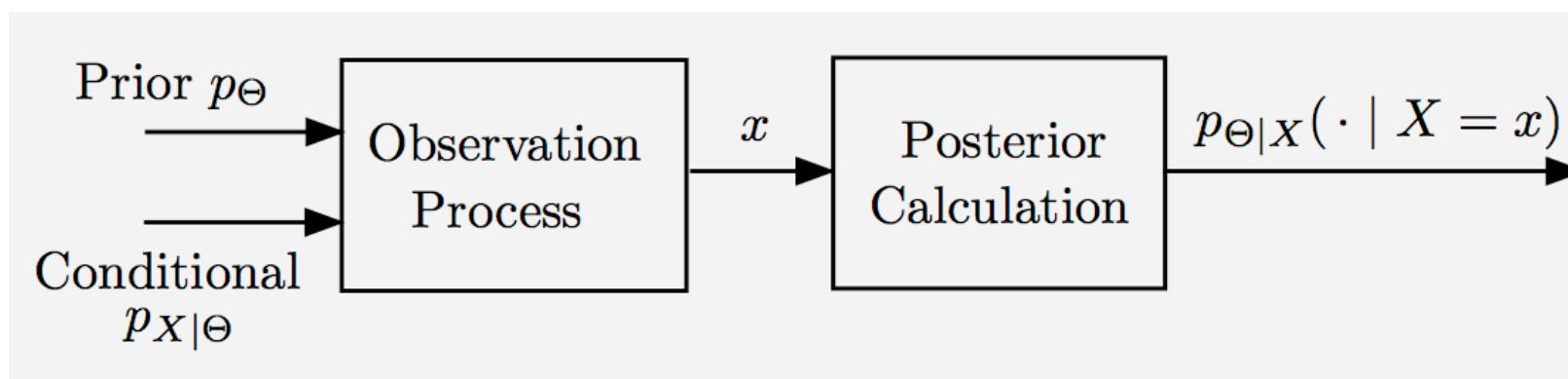
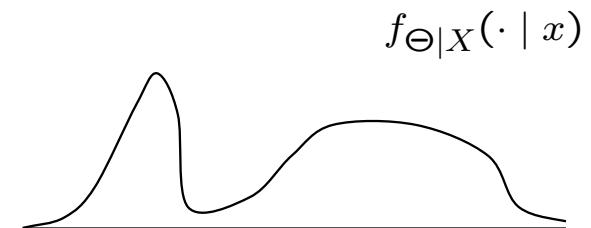
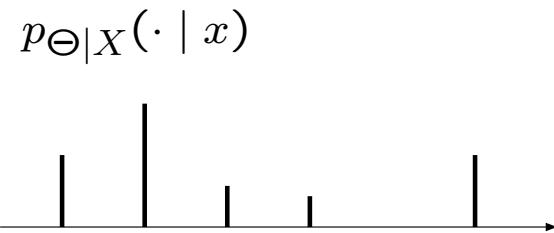
prediction

$$\frac{\mathbb{P}_{\Theta,X}}{X_1 \quad \vdots \quad X_n}$$

inference

The general answer to Bayesian inference problems

$$\frac{\mathbb{P}_{\Theta}}{\Theta?} \quad \mathbb{P}_{X|\Theta}$$



- Use appropriate version of the Bayes rule to find $p_{\Theta|X}(\cdot | X = x)$ or $f_{\Theta|X}(\cdot | X = x)$

$$\mathbb{P}_{\Theta}, \mathbb{P}_{X|\Theta} \rightarrow \mathbb{P}_{X,\Theta} \rightarrow \mathbb{P}_{\Theta|X}$$

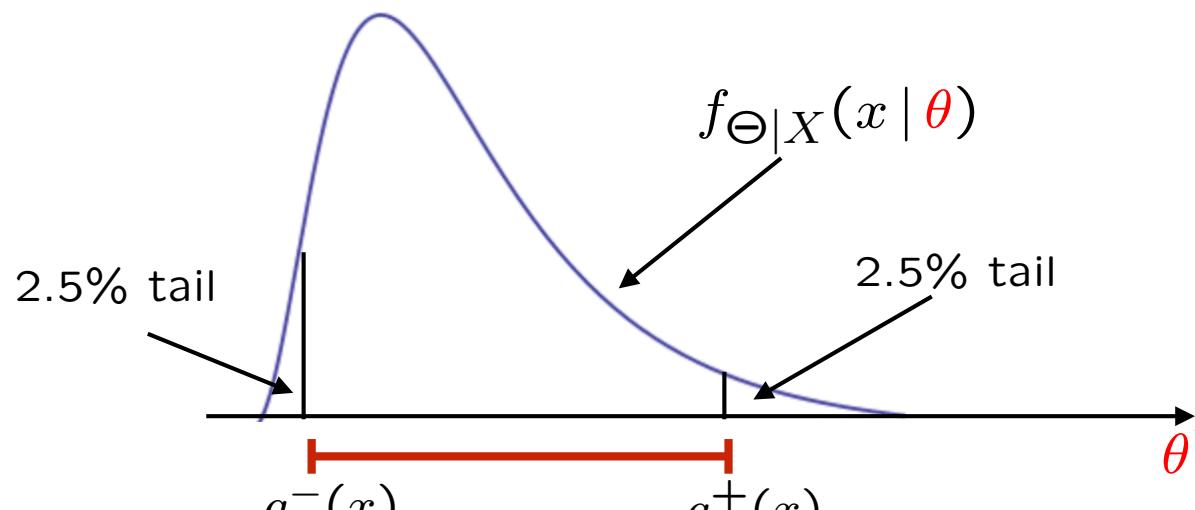
$$\mathbb{P}_{\Theta|X} = \frac{\mathbb{P}_{\Theta} \cdot \mathbb{P}_{X|\Theta}}{\mathbb{P}_X}$$

Four versions of the Bayes' rule

$$\mathbb{P}_{\Theta|X} = \frac{\mathbb{P}_{\Theta} \cdot \mathbb{P}_{X|\Theta}}{\mathbb{P}_X}$$

Variation	Canonical example	Bayes' rule
Θ discrete X discrete	ill/healthy test positive/negative	$p_{\Theta X}(\theta x) = \frac{p_{\Theta}(\theta) \cdot p_{X \Theta}(x \theta)}{p_X(x)}$
Θ discrete X continuous	signal $\Theta=0/1$ $X = \Theta + (\text{independent noise})$	$p_{\Theta X}(\theta x) = \frac{p_{\Theta}(\theta) \cdot f_{X \Theta}(x \theta)}{f_X(x)}$
Θ continuous X discrete	Θ : coin bias $X \sim \text{Binomial}(n, \Theta)$	$f_{\Theta X}(\theta x) = \frac{f_{\Theta}(\theta) \cdot p_{X \Theta}(x \theta)}{p_X(x)}$
Θ continuous X continuous	$X = \Theta + W$ Θ, W : independent, normal	$f_{\Theta X}(\theta x) = \frac{f_{\Theta}(\theta) \cdot f_{X \Theta}(x \theta)}{f_X(x)}$

Bayesian confidence intervals



95% confidence interval

$$\mathbb{P}(g^-(x) \leq \Theta \leq g^+(x) | X = x) = 0.95$$

- Contrast with classical CI:

$$\mathbb{P}^{\theta}(g^-(X) \leq \theta \leq g^+(X)) = 0.95$$

Inferring the unknown bias of a coin

- Coin with bias Θ
 - prior: uniform on $[0, 1]$: $f_{\Theta}(\theta) = 1$, for $0 \leq \theta \leq 1$
- n independent tosses
 - given $\Theta = \theta$, $K \sim \text{Binomial}(n, \theta)$: $p_{K|\Theta}(k | \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$

$$f_{\Theta|K}(\theta | k) = \frac{f_{\Theta}(\theta) p_{K|\Theta}(k | \theta)}{p_K(k)}$$

$$p_K(k) = \int f_{\Theta}(\theta') p_{K|\Theta}(k | \theta') d\theta'$$

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

↑
normalizing constant
dependence on θ

- Animation: <https://av2076.herokuapp.com/bayesian>