

BAYESIAN LEARNING - LECTURE 4

Mattias Villani

**Division of Statistics and Machine Learning
Department of Computer and Information Science
Linköping University**

LECTURE OVERVIEW

► Prediction

- Normal model
- More complex examples

► Decision theory

- The elements of a decision problem
- The Bayesian way
- Point estimation as a decision problem

PREDICTION/FORECASTING

- **Posterior predictive distribution** for future \tilde{y} given observed data y

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta, y) p(\theta|y) d\theta$$

- If $p(\tilde{y}|\theta, y) = p(\tilde{y}|\theta)$ [not true for time series], then

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta) p(\theta|y) d\theta$$

- The **parameter uncertainty** is represented in $p(\tilde{y}|y)$ by **averaging over** $p(\theta|y)$.

PREDICTION - NORMAL DATA, KNOWN VARIANCE

- Under the uniform prior $p(\theta) \propto c$, then

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta)p(\theta|y)d\theta$$

where

$$\begin{aligned}\theta|y &\sim N(\bar{y}, \sigma^2/n) \\ \tilde{y}|\theta &\sim N(\theta, \sigma^2)\end{aligned}$$

PREDICTION - NORMAL DATA, KNOWN VARIANCE

- Under the uniform prior $p(\theta) \propto c$, then

$$p(\tilde{y}|y) = \int_{\theta} p(\tilde{y}|\theta)p(\theta|y)d\theta$$

where

$$\begin{aligned}\theta|y &\sim N(\bar{y}, \sigma^2/n) \\ \tilde{y}|\theta &\sim N(\theta, \sigma^2)\end{aligned}$$

1. Generate a posterior draw of θ ($\theta^{(1)}$) from $N(\bar{y}, \sigma^2/n)$
2. Generate a draw of \tilde{y} ($\tilde{y}^{(1)}$) from $N(\theta^{(1)}, \sigma^2)$ (note the mean)
3. Repeat steps 1 and 2 a large number of times (N) with the result:
 - Sequence of posterior draws: $\theta^{(1)}, \dots, \theta^{(N)}$
 - Sequence of predictive draws: $\tilde{y}^{(1)}, \dots, \tilde{y}^{(N)}$.

PREDICTIVE DISTRIBUTION - NORMAL MODEL AND UNIFORM PRIOR

- ▶ $\theta^{(1)} = \bar{y} + \varepsilon^{(1)}$, where $\varepsilon^{(1)} \sim N(0, \sigma^2/n)$. (Step 1).
- ▶ $\tilde{y}^{(1)} = \theta^{(1)} + v^{(1)}$, where $v^{(1)} \sim N(0, \sigma^2)$. (Step 2).
- ▶ $\tilde{y}^{(1)} = \bar{y} + \varepsilon^{(1)} + v^{(1)}$.
- ▶ $\varepsilon^{(1)}$ and $v^{(1)}$ are independent.
- ▶ The sum of two normal random variables is normal so

$$\begin{aligned} E(\tilde{y}|y) &= \bar{y} \\ V(\tilde{y}|y) &= \frac{\sigma^2}{n} + \sigma^2 = \sigma^2 \left(1 + \frac{1}{n}\right) \end{aligned}$$

$$\tilde{y}|y \sim N \left[\bar{y}, \sigma^2 \left(1 + \frac{1}{n}\right) \right]$$

PREDICTIVE DISTRIBUTION - NORMAL MODEL AND NORMAL PRIOR

- ▶ It is easy to see that the predictive distribution is normal.
- ▶ The mean can be obtained from

$$E_{\tilde{y}|\theta}(\tilde{y}) = \theta$$

and then remove the conditioning on θ by averaging over θ

$$E(\tilde{y}|y) = E_{\theta|y}(\theta) = \mu_n \text{ (Posterior mean of } \theta\text{)}.$$

- ▶ The predictive variance of \tilde{y} (conditional variance formula):

$$\begin{aligned} V(\tilde{y}|y) &= E_{\theta|y}[V_{\tilde{y}|\theta}(\tilde{y})] + V_{\theta|y}[E_{\tilde{y}|\theta}(\tilde{y})] \\ &= E_{\theta|y}(\sigma^2) + V_{\theta|y}(\theta) \\ &= \sigma^2 + \tau_n^2 \\ &= \text{(Population variance + Posterior variance of } \theta\text{)}. \end{aligned}$$

- ▶ In **summary**:

$$\tilde{y}|y \sim N(\mu_n, \sigma^2 + \tau_n^2).$$

BAYESIAN PREDICTION IN MORE COMPLEX MODELS

► Autoregressive process

$$y_t = \phi_1(y_{t-1} - \mu) + \dots + \phi_p(y_{t-p} - \mu) + \varepsilon_t, \quad \varepsilon_t \stackrel{iid}{\sim} N(0, \sigma^2)$$

- Simulate a draw from $p(\phi_1, \phi_2, \dots, \phi_p, \mu, \sigma | y)$
 - Conditional on that draw $\theta^{(1)} = (\phi_1^{(1)}, \phi_2^{(1)}, \dots, \phi_p^{(1)}, \mu^{(1)}, \sigma^{(1)})$, simulate
 - $\tilde{y}_{T+1} \sim p(y_{T+1} | y_T, y_{T-1}, \dots, y_{T-p}, \theta^{(1)})$
 - $\tilde{y}_{T+2} \sim p(y_{T+2} | \tilde{y}_{T+1}, y_T, \dots, y_{T-p}, \theta^{(1)})$
 - and so on.
- Repeat for new θ draws.

BAYESIAN PREDICTION IN MORE COMPLEX MODELS

► Autoregressive process

$$y_t = \phi_1(y_{t-1} - \mu) + \dots + \phi_p(y_{t-p} - \mu) + \varepsilon_t, \quad \varepsilon_t \stackrel{iid}{\sim} N(0, \sigma^2)$$

► Simulate a draw from $p(\phi_1, \phi_2, \dots, \phi_p, \mu, \sigma | y)$

- Conditional on that draw $\theta^{(1)} = (\phi_1^{(1)}, \phi_2^{(1)}, \dots, \phi_p^{(1)}, \mu^{(1)}, \sigma^{(1)})$, simulate
- $\tilde{y}_{T+1} \sim p(y_{T+1} | y_T, y_{T-1}, \dots, y_{T-p}, \theta^{(1)})$
- $\tilde{y}_{T+2} \sim p(y_{T+2} | \tilde{y}_{T+1}, y_T, \dots, y_{T-p}, \theta^{(1)})$
- and so on.

► Repeat for new θ draws.

► Regression trees.

- Uncertainty on which variables to split on, and the split point.
- For given draw of splitting variables and split points, simulate a response. Repeat for many different draws.

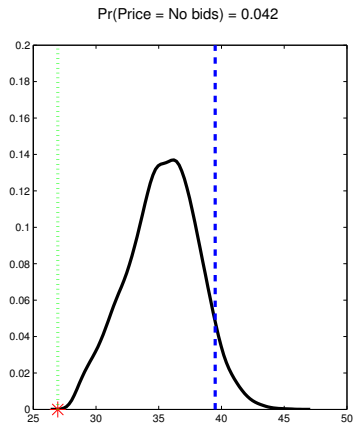
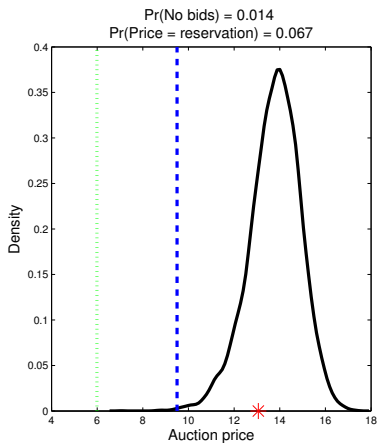
PREDICTING AUCTION PRICES ON EBAY

- ▶ Problem: Predicting the auctioned price in eBay coin auctions.
- ▶ Data: Bid from 1000 auctions on eBay.
 - ▶ The highest bid is not observed.
 - ▶ The lowest bids are also not observed because of the seller's reservation price.
- ▶ Covariates: auction-specific, e.g. Book value from catalog, seller's reservation price, quality of sold object, rating of seller, powerseller, verified seller ID etc
- ▶ Buyers are strategic. Their bids does not fully reflect their valuation. Game theory. Very complicated likelihood.

SIMULATING AUCTION PRICES ON EBAY, CONT.

- ▶ A draw from the **posterior predictive distribution** of an auction's price:
 1. Simulate a draw $\theta^{(1)}$ from the posterior of the model parameters θ (using MCMC)
 2. Simulate the number of bidders conditional on θ (Poisson process)
 3. Simulate the bidders' valuations.
 4. Simulate a complete auction bid sequence, $\mathbf{b}^{(1)}$, conditional on the valuations and $\theta = \theta^{(1)}$.
 5. For the bid sequence $\mathbf{b}^{(1)}$, return the next to largest bid (eBay's proxy bidding system).

PREDICTING AUCTION PRICES ON EBAY, CONT.



DECISION THEORY

- ▶ Let θ be an **unknown quantity**. **State of nature**. Examples: Future inflation, Global temperature, Disease.
- ▶ Let $a \in \mathcal{A}$ be an **action**. Ex: Interest rate, Energy tax, Surgery.
- ▶ Choosing action a when state of nature turns out to be θ gives **utility**

$$U(a, \theta)$$

- ▶ Alternatively **loss** $L(a, \theta) = -U(a, \theta)$.

- ▶ Loss table:

	θ_1	θ_2
a_1	$L(a_1, \theta_1)$	$L(a_1, \theta_2)$
a_2	$L(a_2, \theta_1)$	$L(a_2, \theta_2)$

- ▶ Example:

	Rainy	Sunny
Umbrella	20	10
No umbrella	50	0

DECISION THEORY, CONT.

- ▶ Example **loss functions** when both a and θ are continuous:

- ▶ **Linear**: $L(a, \theta) = |a - \theta|$
- ▶ **Quadratic**: $L(a, \theta) = (a - \theta)^2$
- ▶ **Lin-Lin**:

$$L(a, \theta) = \begin{cases} c_1 \cdot |a - \theta| & \text{if } a \leq \theta \\ c_2 \cdot |a - \theta| & \text{if } a > \theta \end{cases}$$

- ▶ Example:

- ▶ θ is the number of items demanded of a product
- ▶ a is the number of items in stock
- ▶ Utility

$$U(a, \theta) = \begin{cases} p \cdot \theta - c_1(a - \theta) & \text{if } a > \theta \text{ [too much stock]} \\ p \cdot a - c_2(\theta - a)^2 & \text{if } a \leq \theta \text{ [too little stock]} \end{cases}$$

OPTIMAL DECISION

- ▶ Ad hoc decision rules:
 - ▶ *Minimax*. Choose the decision that minimizes the maximum loss.
 - ▶ *Minimax-regret* ... bla bla bla ...
- ▶ Bayesian **theory**: Just maximize the **posterior expected utility**:

$$a_{\text{bayes}} = \operatorname{argmax}_{a \in \mathcal{A}} E_{p(\theta|y)}[U(a, \theta)],$$

where $E_{p(\theta|y)}$ denotes the posterior expectation.

- ▶ Using simulated draws $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(N)}$ from $p(\theta|y)$:

$$E_{p(\theta|y)}[U(a, \theta)] \approx N^{-1} \sum_{i=1}^N U(a, \theta^{(i)})$$

- ▶ **Separation principle**:
 1. First obtain $p(\theta|y)$
 2. then form $U(a, \theta)$ and finally
 3. choose a that maximizes $E_{p(\theta|y)}[U(a, \theta)]$.

CHOOSING A POINT ESTIMATE IS A DECISION

- ▶ Choosing a **point estimator** is a decision problem.
- ▶ Which to choose: posterior median, mean or mode?
- ▶ It depends on your loss function:
 - ▶ **Linear loss** → Posterior median is optimal
 - ▶ **Quadratic loss** → Posterior mean is optimal
 - ▶ **Lin-Lin loss** → $c_2 / (c_1 + c_2)$ quantile of the posterior is optimal
 - ▶ **Zero-one loss** → Posterior mode is optimal