# Bayesian Learning
## Lecture 4 - Regression, Prediction and Decisions

**Mattias Villani** 🧔

Department of Statistics
Stockholm University

Department of Computer and Information Science
Linköping University

🌐 mattiasvillani.com        🐦 @matvil        ⭘ mattiasvillani

# Lecture overview

- **Normal model** with conjugate prior

- The **linear regression** model

- **Prediction**

- **Decision making**

# Linear regression

■ The linear regression model in **matrix form**

$$\underset{(n\times 1)}{y} = \underset{(n\times k)(k\times 1)}{X\beta} + \underset{(n\times 1)}{\varepsilon}$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \ \beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix}, \ \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$X = \begin{pmatrix} x_1' \\ \vdots \\ x_n' \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1k} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nk} \end{pmatrix}$$

■ Usually $x_{i1} = 1$, for all $i$. $\beta_1$ is the intercept.

■ **Likelihood**

$$y|\beta, \sigma^2, X \sim N(X\beta, \sigma^2 I_n)$$

# Linear regression - uniform prior

- Standard non-informative prior: uniform on $(\beta, \log \sigma^2)$

$$p(\beta, \sigma^2) \propto \sigma^{-2}$$

- Joint posterior of $\beta$ and $\sigma^2$:

$$
\begin{aligned}
\beta | \sigma^2, y &\sim N\left[\hat{\beta}, \sigma^2 (X'X)^{-1}\right] \\
\sigma^2 | y &\sim \text{Inv-}\chi^2(n-k, s^2)
\end{aligned}
$$

where $\hat{\beta} = (X'X)^{-1} X'y$ and $s^2 = \frac{1}{n-k}(y - X\hat{\beta})'(y - X\hat{\beta})$.

- Simulate from the joint posterior by simulating from
  - $p(\sigma^2 | y)$
  - $p(\beta | \sigma^2, y)$

- Marginal posterior of $\beta$:

$$\beta | y \sim t_{n-k}\left[\hat{\beta}, s^2 (X'X)^{-1}\right]$$

# Linear regression - conjugate prior

- **Joint prior** for $\beta$ and $\sigma^2$

$$\beta | \sigma^2 \sim N\left(\mu_0, \sigma^2 \Omega_0^{-1}\right)$$
$$\sigma^2 \sim Inv - \chi^2\left(\nu_0, \sigma_0^2\right)$$

- **Posterior**

$$\beta | \sigma^2, \mathbf{y} \sim N\left[\mu_n, \sigma^2 \Omega_n^{-1}\right]$$
$$\sigma^2 | \mathbf{y} \sim \text{Inv} - \chi^2\left(\nu_n, \sigma_n^2\right)$$

$$\mu_n = \left(\mathsf{X}^\top \mathsf{X} + \Omega_0\right)^{-1}\left(\mathsf{X}^\top \mathsf{X} \hat{\beta} + \Omega_0 \mu_0\right)$$
$$\Omega_n = \mathsf{X}^\top \mathsf{X} + \Omega_0$$
$$\nu_n = \nu_0 + n$$
$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + \left(\mathbf{y}^\top \mathbf{y} + \mu_0^\top \Omega_0 \mu_0 - \mu_n^\top \Omega_n \mu_n\right)$$
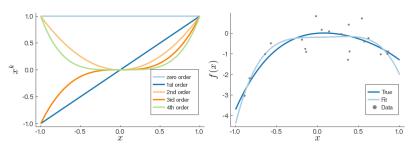
# Polynomial regression

■ **Polynomial regression**

$$f(x_i) = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + ... + \beta_k x_i^k.$$

$$y = X_P \beta + \varepsilon,$$

where

$$X_P = (1, x, x^2, ..., x^k).$$



■ Priors for **regularization** (ridge, lasso etc) in Lecture 6.

# Prediction/Forecasting

- **Posterior predictive density** for future $\tilde{y}$ given observed $\boldsymbol{y} = (y_1, \ldots, y_n)$

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

- IID data:

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

- **Parameter uncertainty** in $p(\tilde{y}|\boldsymbol{y})$ by **averaging over** $p(\theta|\boldsymbol{y})$.

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

$$p(\tilde{y}|\boldsymbol{y}) = \int_{\theta} p(\tilde{y}|\theta)p(\theta|\boldsymbol{y})d\theta$$
$$\theta|\boldsymbol{y} \sim N(\bar{y}, \sigma^2/n)$$
$$\tilde{y}|\theta \sim N(\theta, \sigma^2)$$

Simulation algorithm:

1. Generate a **posterior draw** of $\theta$ ($\theta^{(1)}$) from $N(\bar{y}, \sigma^2/n)$
2. Generate a **predictive draw** of $\tilde{y}$ ($\tilde{y}^{(1)}$) from $N(\theta^{(1)}, \sigma^2)$
3. Repeat Steps 1 and 2 $N$ times to output:
   - ▶ Sequence of posterior draws: $\theta^{(1)}, ...., \theta^{(N)}$
   - ▶ Sequence of predictive draws: $\tilde{y}^{(1)}, ..., \tilde{y}^{(N)}$.

# Predictive distribution - Normal model

- $\theta^{(1)} = \bar{y} + \varepsilon^{(1)}$, 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

$$E(\tilde{y}|\boldsymbol{y}) = \bar{y}$$

$$V(\tilde{y}|\boldsymbol{y}) = \frac{\sigma^2}{n} + \sigma^2 = \sigma^2 \left(1 + \frac{1}{n}\right)$$

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

# Iteration laws

- Expectation with respect to what? Explicit:

$$\mathbb{E}_{\theta|\boldsymbol{y}}(\theta) \equiv \int \theta p(\theta|\boldsymbol{y})d\theta$$

- Law of iterated expectation and Law of total variance.

---

**Iteration laws**

Law of iterated expectation:

$$\mathbb{E}_X(X) = \mathbb{E}_Y\big(\mathbb{E}_{X|Y}(X)\big)$$

Law of total variance:

$$\mathbb{V}_X(X) = \mathbb{E}_Y\big(\mathbb{V}_{X|Y}(X)\big)$$
$$+\mathbb{V}_Y\big(\mathbb{E}_{X|Y}(X)\big)$$

---

**Iteration laws for Bayes**

Marginal posterior mean:

$$\mathbb{E}_{\theta_1|\mathbf{y}}(\theta_1) = \mathbb{E}_{\theta_2|\mathbf{y}}\big(\mathbb{E}_{\theta_1|\theta_2,\mathbf{y}}(\theta_1)\big)$$

Marginal posterior variance:

$$\mathbb{V}_{\theta_1}(\theta_1) = \mathbb{E}_{\theta_2|\mathbf{y}}\big(\mathbb{V}_{\theta_1|\theta_2,\mathbf{y}}(\theta_1)\big)$$
$$+\mathbb{V}_{\theta_2|\mathbf{y}}\big(\mathbb{E}_{\theta_1|\theta_2,\mathbf{y}}(\theta_1)\big)$$

# Predictive distribution - Normal model and prior

■ Predictive distribution still normal (sum of normals is normal).

■ Predictive mean conditional on $\theta$ is trivial:

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

■ "Remove the conditioning" on $\theta$ by averaging over posterior:

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

■ The predictive variance of $\tilde{y}$ by law of total variance

$$\begin{aligned} V(\tilde{y}|\boldsymbol{y}) &= E_{\theta|\boldsymbol{y}}[V_{\tilde{y}|\theta}(\tilde{y})] + V_{\theta|\boldsymbol{y}}[E_{\tilde{y}|\theta}(\tilde{y})] \\ &= E_{\theta|\boldsymbol{y}}(\sigma^2) + V_{\theta|\boldsymbol{y}}(\theta) \\ &= \sigma^2 + \tau_n^2 \end{aligned}$$
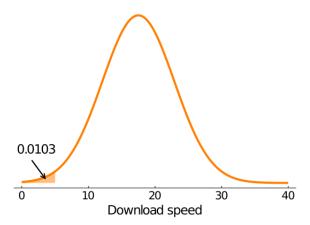
■ So, predictive distribution is

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

# Predictive distribution – Internet speed data
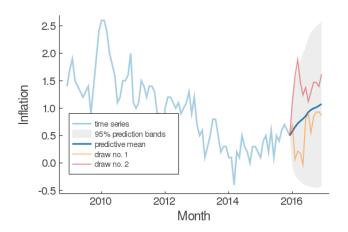
- My Netflix starts to buffer at speeds $< 5$Mbit.

# Bayesian prediction for time series

■ **Autoregressive process**

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

---

**Predictive distribution - AR process.**

**Input:** time series $\mathbf{y}_{1:T} = (y_1, \ldots, y_T)$
number of predictive draws $m$.
forecast horizon $h$.

**for** $i$ in $1{:}m$ **do**

$\quad \mu, \phi_1, \ldots, \phi_p, \sigma \leftarrow \text{RPOSTERIORAR}(\mathbf{y}_{1:T}, \text{PriorSettings})$

$\quad \varepsilon_{T+1} \leftarrow \text{RNORM}(0, \sigma)$
$\quad \tilde{y}_{T+1} \leftarrow \mu + \phi_1(y_T - \mu) + \ldots + \phi_p(y_{T+1-p} - \mu) + \varepsilon_{T+1}$
$\quad \varepsilon_{T+2} \leftarrow \text{RNORM}(0, \sigma)$
$\quad \tilde{y}_{T+2} \leftarrow \mu + \phi_1(\tilde{y}_{T+1} - \mu) + \ldots + \phi_p(y_{T+2-p} - \mu) + \varepsilon_{T+2}$
$\quad \qquad \vdots$
$\quad \varepsilon_{T+h} \leftarrow \text{RNORM}(0, \sigma)$
$\quad \tilde{y}_{T+h} \leftarrow \mu + \phi_1(\tilde{y}_{T+h-1} - \mu) + \ldots + \phi_p(\tilde{y}_{T+h-p} - \mu) + \varepsilon_{T+h}$

**end**

**Output:** $m$ draws from the joint predictive density:
$$p(\tilde{y}_{T+1}, \ldots, \tilde{y}_{T+h} | \mathbf{y}_{1:T}).$$

# Bayesian prediction of Swedish inflation

# Predicting auction prices on eBay

- **Problem:** **Predicting the final price** in eBay coin auctions.

- **Data:** Bid from 1000 auctions on eBay.
  The highest bid is not observed (eBay proxy bidding).

- **Covariates** are auction-specific:
  - catalog value
  - seller's **reservation price**
  - quality
  - rating of seller etc

- Buyers are **strategic**.
  - Bid $\neq$ **valuation**.
  - **Bid function**, $b = \text{BidFunction}(v)$, from **Game theory**.
  - Very complicated likelihood.

# Simulating auction prices on eBay
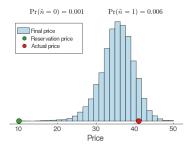
**Predictive distribution - auction price.**

**Input:** training auction bids $\mathbf{Y}$
training auction covariates $\mathbf{X}$.
test auction covariates $\tilde{\mathbf{x}}$.
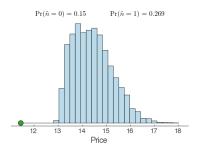number of predictive draws $m$.

**for** $i$ in $1\!:\!m$ **do**

$\mu, \sigma, \lambda \leftarrow \text{RPOSTAUCTION}(\mathbf{Y}, \mathbf{X}, \tilde{\mathbf{x}}, \text{Prior})$ # parameters

$\tilde{n} \leftarrow \text{RPOIS}(\lambda(\tilde{\mathbf{x}}))$ # number of bidders in test auction

$\tilde{\mathbf{v}}_{1:\tilde{n}} \leftarrow \text{RNORM}(\mu(\tilde{\mathbf{x}}), \sigma(\tilde{\mathbf{x}}))$ # valuations for all $\tilde{n}$ bidders

$b_{1:\tilde{n}} \leftarrow \text{BIDFUNCTION}(\tilde{\mathbf{v}}_{1:\tilde{n}}, \tilde{n}, \mu(\tilde{\mathbf{x}}), \sigma(\tilde{\mathbf{x}}))$ # bids

$\tilde{p} \leftarrow \text{SECONDLARGEST}(b_{1:\tilde{n}})$ # final price

**end**

**Output:** $m$ predictive draws of the final price $\tilde{p}$ for an
auction with covariates $\tilde{\mathbf{x}}$.

# Predicting auction prices on eBay

# Decision problems

- Let $\theta$ be an **unknown quantity**. **State of nature**.
  - ▶ Future inflation
  - ▶ Global temperature
  - ▶ Disease.

- Let $a \in \mathcal{A}$ be an **action**.
  - ▶ Interest rate
  - ▶ Energy tax
  - ▶ Surgery.

- Choosing action $a$ when state of nature is $\theta$ gives **utility**

$$U(a, \theta)$$

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

■ **Decision table**

|  | $\theta_1$ | $\theta_2$ | $\cdots$ | $\theta_K$ |
|---|---|---|---|---|
| $a_1$ | $u(a_1, \theta_1)$ | $u(a_1, \theta_2)$ | $\cdots$ | $u(a_1, \theta_K)$ |
| $a_2$ | $u(a_2, \theta_1)$ | $u(a_2, \theta_2)$ | $\cdots$ | $u(a_2, \theta_K)$ |
| $\vdots$ | $\vdots$ | $\vdots$ |  | $\vdots$ |
| $a_J$ | $u(a_J, \theta_1)$ | $u(a_J, \theta_2)$ | $\cdots$ | $u(a_J, \theta_K)$ |

■ The eternal umbrella decision:

|  | Rain | Sun |
|---|---|---|
| No umbrella | $-50$ | 50 |
| Umbrella | 10 | 30 |

# Decision Theory

- Example loss functions when both $a$ and $\theta$ 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:
  - $\theta$ 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 decisions

- Ad hoc decision rules:
    - *Minimax*. Minimizes the maximum loss.
    - *Minimax-regret*
    - ... 😴
- Bayesian theory: maximize posterior expected utility 😍

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

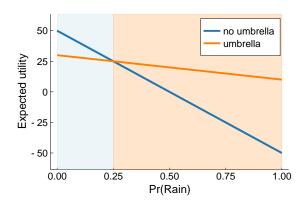where $E_{p(\theta|y)}$ denotes the posterior expectation.
- Using simulated draws $\theta^{(1)}, \theta^{(2)}, ..., \theta^{(N)}$ from $p(\theta|y)$ :

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

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

# The umbrella decision

|              | Rain | Sun |
|--------------|------|-----|
| No umbrella  | −50  | 50  |
| Umbrella     | 10   | 30  |

# 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
  - ▶ **Quadratic loss** → Posterior mean
  - ▶ **Zero-one loss** → Posterior mode
  - ▶ **Lin-Lin loss** → $c_1/(c_1 + c_2)$ quantile of the posterior

# The umbrella decision