# Bayesian Learning Lecture 11 - Course summary

#### Mattias Villani

Department of Statistics Stockholm University











### **Overview**

■ Wait, what did we actually do?

## **Bayesics**

- Subjective probability
- Bayesian Learning: Bayes theorem

$$p(\theta|\mathbf{x}) = \frac{p(\mathbf{x}|\theta)p(\theta)}{p(\mathbf{x})} \propto p(\mathbf{x}|\theta)p(\theta)$$

■ The proportional constant is actually the marginal likelihood

$$p(\mathbf{x}) = \int p(\mathbf{x}|\theta)p(\theta)d\theta$$

used for model comparison.

## Simple models - conjugate priors

- iid Bernoulli model Beta prior
- iid Poisson model Gamma prior
- iid Exponential model Gamma prior
- iid Normal known variance Normal prior

## Multi-parameter models

Marginalization

$$p( heta_1| extbf{ extit{x}}) = \int p( heta_1, heta_2| extbf{ extit{x}}) extit{d} heta_2$$

iid normal with unknown variance  $N(\theta, \sigma^2)$ . Marginal posterior for  $\theta$  is student-t

$$\theta | \mathbf{x} \sim t \left( \mu_{\mathbf{n}}, \frac{\sigma_{\mathbf{n}}^2}{\kappa_{\mathbf{n}}}, \nu_{\mathbf{n}} \right)$$

- Multinomial model Dirichlet prior
- Dirichlet distribution on unit simplex  $\sum_{k=1}^{K} \theta_k = 1$ .

## **Linear Gaussian regression**

Linear regression

$$y_i = \mathbf{x}_i^{\top} \boldsymbol{\beta} + \varepsilon_i, \qquad \varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

Conjugate prior

$$oldsymbol{eta} | oldsymbol{\sigma}^2 \sim oldsymbol{N} \left( \mu_0, \sigma^2 \Omega_0^{-1} 
ight) \ \sigma^2 \sim \operatorname{Inv} - \chi^2 \left( \nu_0, \sigma_0^2 
ight)$$

Posterior

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

L2-Regularization prior (Ridge)

$$m{eta}|\sigma^2 \sim m{N}igg(m{0}, rac{\sigma^2}{\lambda}m{I}_pigg)$$

L1-Regularization prior (Lasso)

$$oldsymbol{eta} | oldsymbol{eta}^2 \sim ext{Laplace} \left( oldsymbol{0}, rac{\sigma^2}{\lambda} oldsymbol{I}_{\! p} 
ight)$$

## Logistic regression and Beyond

Logistic regression

$$\Pr(y_i = 1 | \mathbf{x}_i) = \frac{e^{\mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta}}}{1 + e^{\mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta}}}$$

- Posterior distribution  $p(\beta|\mathbf{y},\mathbf{X})$  is intractable. Solutions:
  - Normal approximation
  - Posterior sampling using MH/HMC
- Linear Gaussian regression is modeling a conditional density

$$y_i | \boldsymbol{x}_i \overset{\text{ind}}{\sim} N(\mu_i, \sigma^2) \text{ where } \mu_i = \boldsymbol{x}_i^{\top} \boldsymbol{\beta}$$

Paves the way for Poisson regression:

$$y_i | \mathbf{x}_i \overset{\text{ind}}{\sim} \operatorname{Pois}(\mu_i) \text{ where } \mu_i = e^{\mathbf{x}_i^{\top} \boldsymbol{\beta}}$$

Exponential regression

$$y_i | \mathbf{x}_i \overset{\mathrm{ind}}{\sim} \mathrm{Expon}(\theta_i) \ \mathsf{where} \ \theta_i = \frac{1}{e^{\mathbf{x}_i^{\mathsf{T}} \boldsymbol{\beta}}}$$

so that 
$$\mathbb{E}(y_i|\mathbf{x}_i) = \mu_i = e^{\mathbf{x}_i^{\top}\boldsymbol{\beta}}$$
.

#### **Priors**

- **Expert** elicitation
- Other data sources
- Non-informative priors (prior sample size). Bernoulli model:

$$\theta | \mathbf{x} \sim \text{Beta}(\alpha + \mathbf{s}, \beta + \mathbf{f})$$

- Jeffreys' prior  $p(\theta) = |I(\theta)|^{1/2}$
- Hierarchical priors. L2-regularization

$$p(\beta, \sigma^2, \lambda) = p(\beta | \sigma^2 \lambda) p(\sigma^2) p(\lambda)$$

Regularization and smoothness priors

#### **Prediction**

(Posterior) predictive distribution for iid data:

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

- $ightharpoonup p(\tilde{y}|\theta)$  is the model density for a new observation  $\tilde{y}$  given  $\theta$
- $ightharpoonup p( heta|m{y})$  is the posterior density for heta given the training data  $m{y}$
- With no training data: the prior predictive distribution

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

where we integrate of the parameters using the **prior**  $p(\theta)$ .

- Use prior predictive density to set the prior hyperparameters when the expert expresses prior beliefs about observable data:
  - $\triangleright$   $\mathbb{E}(\tilde{y})$
  - $ightharpoonup \Pr(\tilde{y} > c)$
- Interactive example with Gamma prior for Poisson model.

#### **Decisions**

- Need to make a decision/action  $a \in A$  when state of nature  $\theta$  is partially unknown.
- The utility function describes the consequence of taking action a when state of nature is  $\theta$

$$U(a, \theta)$$

- The utility function is subjective (also for non-Bayesians) and is determined by the decision maker.
- Bayesian theory gives a single universal decision rule: choose the action that maximizes (posterior) expected utility:

$$\mathrm{EU}(a) = \int U(a, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\boldsymbol{y}) d\boldsymbol{\theta}$$

where y is some training data.

## **Posterior approximation**

- Three ways to approximate the posterior:
  - ► Normal approximation
  - Posterior sampling (Gibbs, MH, HMC)
  - Numerical integration
- Normal approximation is fast, but is approximate unless we have a lot of data (infinite in theory).
- Posterior sampling will converge to true posterior if simulated long enough, but is approximate in finite time.
- Numerical integration is only a good option when the number of parameters is small.

## Model evaluation and comparison

- Three main ways to compare models
  - ► Posterior model probabilities (marginal likelihood)
  - ▶ Log Predictive Density Score (LPDS/elpd)
  - ► Leave-one-out (LOO) LPDS/elpd
- Evaluate the predictive density performance on test data

$$p(\mathbf{y}_{ ext{test}}|\mathbf{y}_{ ext{train}}) = \int p(\mathbf{y}_{ ext{test}}|\mathbf{\theta}, \mathbf{y}_{ ext{train}}) p(\mathbf{\theta}|\mathbf{y}_{ ext{train}}) d\mathbf{\theta}$$

- test data y<sub>test</sub>
- lacktriangle posterior  $p(m{ heta}|m{y}_{\mathrm{train}})$  based on some training data.
- Often evaluated on the log scale:  $\log p(\mathbf{y}_{\text{test}}|\mathbf{y}_{\text{train}})$ .
- Parameter uncertainty is included in the model comparison by integrating with respect to the posterior  $p(\theta|\mathbf{y}_{\text{train}})$ .

## Model evaluation and comparison

#### Marginal likelihood:

- ightharpoonup no  $y_{\text{train}}$  and  $y_{\text{test}}$  is all the data.
- ▶ Sensitive to the prior on the parameters  $p(\theta)$ .

#### ■ Log Predictive Density Score

- $ightharpoonup y_{\text{train}}$  and  $y_{\text{test}}$  is some partition of the data.
- K-fold cross-validation.

#### Leave-one-out (LOO)

- $\mathbf{y}_{\text{test}} = y_i \text{ and } \mathbf{y}_{\text{train}} = \mathbf{y}_{-i} \text{ (all data except obs } i\text{)}$
- n-fold cross-validation.
- ▶ Time-consuming: posterior sampling from *n* posteriors

$$p(\boldsymbol{\theta}|\mathbf{y}_{-i})$$
 for  $i=1,\ldots,n$ 

▶ The **loo** package uses a importance sampling trick for speed.