

## HT2015: SC4 Statistical Data Mining and Machine Learning

Dino Sejdinovic  
Department of Statistics  
Oxford

<http://www.stats.ox.ac.uk/~sejdinov/sdmml.html>

## Bayesian Learning

### Maximum Likelihood Principle

- A generative model for training data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  given a parameter vector  $\theta$ :

$$y_i \sim (\pi_1, \dots, \pi_K), \quad x|y_i \sim g_{y_i}(x) = p(x|\phi_{y_i})$$

- $k$ -th class conditional density assumed to have a parametric form for  $g_k(x) = p(x|\phi_k)$  and all parameters are given by  $\theta = (\pi_1, \dots, \pi_K; \phi_1, \dots, \phi_K)$
- Generative process defines the **likelihood function**: the joint distribution of all the observed data  $p(\mathcal{D}|\theta)$  given a parameter vector  $\theta$ .
- Process of generative learning consists of computing the MLE  $\hat{\theta}$  of  $\theta$  based on  $\mathcal{D}$ :

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmax}} p(\mathcal{D}|\theta)$$

- We then use a plug-in approach to perform classification

$$f_{\hat{\theta}}(x) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \mathbb{P}_{\hat{\theta}}(Y = k|X = x) = \underset{k \in \{1, \dots, K\}}{\operatorname{argmax}} \frac{\hat{\pi}_k p(x|\hat{\phi}_k)}{\sum_{j=1}^K \hat{\pi}_j p(x|\hat{\phi}_j)}$$

### The Bayesian Learning Framework

- Being Bayesian: **treat parameter vector  $\theta$  as a random variable**: process of learning is then **computation of the posterior distribution**  $p(\theta|\mathcal{D})$ .
- In addition to the likelihood  $p(\mathcal{D}|\theta)$  need to specify a **prior distribution**  $p(\theta)$ .

- Posterior distribution is then given by the **Bayes Theorem**:

$$p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{p(\mathcal{D})}$$

- **Likelihood**:  $p(\mathcal{D}|\theta)$
- **Posterior**:  $p(\theta|\mathcal{D})$
- **Prior**:  $p(\theta)$
- **Marginal likelihood**:  $p(\mathcal{D}) = \int_{\Theta} p(\mathcal{D}|\theta)p(\theta)d\theta$

- Summarizing the posterior:

- **Posterior mode**:  $\hat{\theta}^{\text{MAP}} = \operatorname{argmax}_{\theta \in \Theta} p(\theta|\mathcal{D})$  (maximum a posteriori).
- **Posterior mean**:  $\hat{\theta}^{\text{mean}} = \mathbb{E}[\theta|\mathcal{D}]$ .
- **Posterior variance**:  $\operatorname{Var}[\theta|\mathcal{D}]$ .

## Simple Example: Coin Tosses

- A simple example: We have a coin with probability  $\phi$  of coming up heads. Model coin tosses as i.i.d. Bernoullis,  $1 = \text{head}$ ,  $0 = \text{tail}$ .
- Estimate  $\phi$  given a dataset  $\mathcal{D} = \{x_i\}_{i=1}^n$  of tosses.

$$p(\mathcal{D}|\phi) = \phi^{n_1} (1 - \phi)^{n_0}$$

with  $n_j = \sum_{i=1}^n \mathbb{1}(x_i = j)$ .

- Maximum Likelihood estimate:

$$\hat{\phi}^{\text{ML}} = \frac{n_1}{n}$$

- Bayesian approach: treat the unknown parameter  $\phi$  as a random variable. Simple prior:  $\phi \sim \text{Uniform}[0, 1]$ , i.e.,  $p(\phi) = 1$  for  $\phi \in [0, 1]$ . Posterior distribution:

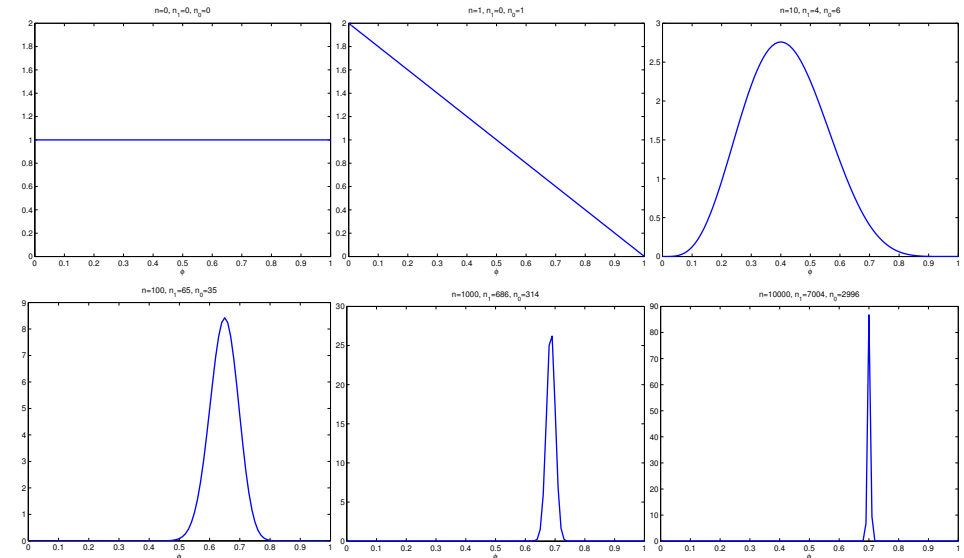
$$p(\phi|\mathcal{D}) = \frac{p(\mathcal{D}|\phi)p(\phi)}{p(\mathcal{D})} = \frac{\phi^{n_1}(1-\phi)^{n_0} \cdot 1}{p(\mathcal{D})}, \quad p(\mathcal{D}) = \int_0^1 \phi^{n_1}(1-\phi)^{n_0} d\phi = \frac{(n+1)!}{n_1!n_0!}$$

Posterior is a  $\text{Beta}(n_1 + 1, n_0 + 1)$  distribution:  $\hat{\phi}^{\text{mean}} = \frac{n_1 + 1}{n + 2}$ .

## Simple Example: Coin Tosses

- All Bayesian reasoning is based on the posterior distribution.
  - Posterior mode:  $\hat{\phi}^{\text{MAP}} = \frac{n_1}{n}$
  - Posterior mean:  $\hat{\phi}^{\text{mean}} = \frac{n_1 + 1}{n + 2}$
  - Posterior variance:  $\text{Var}[\phi|\mathcal{D}] = \frac{1}{n+3} \hat{\phi}^{\text{mean}} (1 - \hat{\phi}^{\text{mean}})$
  - $(1 - \alpha)$ -credible regions:  $(l, r) \subset [0, 1]$  s.t.  $\int_l^r p(\theta|\mathcal{D}) d\theta = 1 - \alpha$ .
- Consistency: Assuming that the true parameter value  $\phi^*$  is given a non-zero density under the prior, the posterior distribution concentrates around the true value as  $n \rightarrow \infty$ .
- Rate of convergence?

## Simple Example: Coin Tosses



Posterior becomes behaves like the ML estimate as dataset grows and is peaked at true value  $\phi^* = .7$ .

## Simple Example: Coin Tosses

- The **posterior predictive distribution** is the conditional distribution of  $x_{n+1}$  given  $\mathcal{D} = \{x_i\}_{i=1}^n$ :

$$\begin{aligned} p(x_{n+1}|\mathcal{D}) &= \int_0^1 p(x_{n+1}|\phi, \mathcal{D}) p(\phi|\mathcal{D}) d\phi \\ &= \int_0^1 p(x_{n+1}|\phi) p(\phi|\mathcal{D}) d\phi \\ &= (\hat{\phi}^{\text{mean}})^{x_{n+1}} (1 - \hat{\phi}^{\text{mean}})^{1-x_{n+1}} \end{aligned}$$

- We predict on new data by **averaging** the predictive distribution over the posterior. Accounts for uncertainty about  $\phi$ .

## Simple Example: Coin Tosses

- In this example, the posterior distribution has a known analytic form and is in the same Beta family as the prior:  $\text{Uniform}[0, 1] \equiv \text{Beta}(1, 1)$ .
- An example of a **conjugate prior**.
- A Beta distribution  $\text{Beta}(a, b)$  with parameters  $a, b > 0$  is an exponential family distribution with density

$$p(\phi|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \phi^{a-1} (1-\phi)^{b-1}$$

where  $\Gamma(t) = \int_0^\infty u^{t-1} e^{-u} du$  is the gamma function.

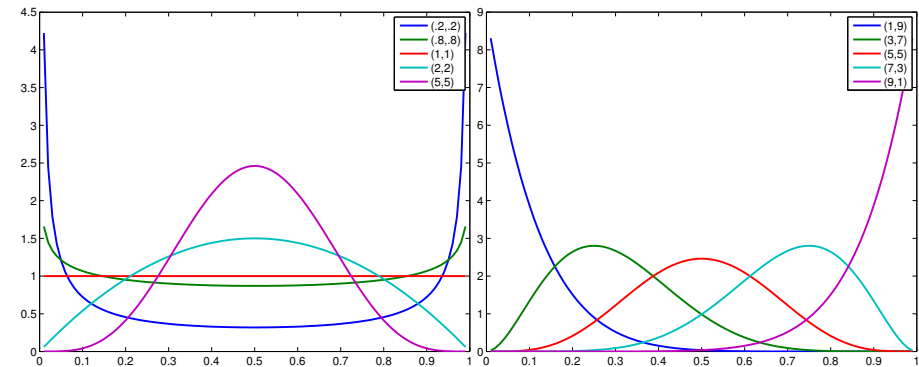
- If the prior is  $\phi \sim \text{Beta}(a, b)$ , then the posterior distribution is

$$p(\phi|\mathcal{D}, a, b) \propto \phi^{a+n_1-1} (1-\phi)^{b+n_0-1}$$

so is  $\text{Beta}(a+n_1, b+n_0)$ .

- Hyperparameters  $a$  and  $b$  are **pseudo-counts**, an imaginary initial sample that reflects our prior beliefs about  $\phi$ .

## Beta Distributions



## Bayesian Inference on the Categorical Distribution

- Suppose we observe  $\mathcal{D} = \{y_i\}_{i=1}^n$  with  $y_i \in \{1, \dots, K\}$ , and model them as i.i.d. with pmf  $\pi = (\pi_1, \dots, \pi_K)$ :

$$p(\mathcal{D}|\pi) = \prod_{i=1}^n \pi_{y_i} = \prod_{k=1}^K \pi_k^{n_k}$$

with  $n_k = \sum_{i=1}^n \mathbb{1}(y_i = k)$  and  $\pi_k > 0$ ,  $\sum_{k=1}^K \pi_k = 1$ .

- The conjugate prior on  $\pi$  is the Dirichlet distribution  $\text{Dir}(\alpha_1, \dots, \alpha_K)$  with parameters  $\alpha_k > 0$ , and density

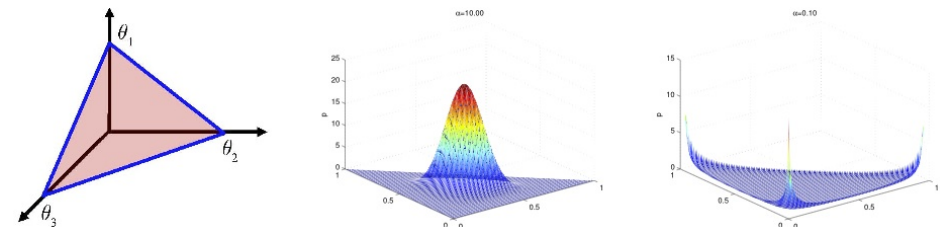
$$p(\pi) = \frac{\Gamma(\sum_{k=1}^K \alpha_k)}{\prod_{k=1}^K \Gamma(\alpha_k)} \prod_{k=1}^K \pi_k^{\alpha_k-1}$$

on the probability simplex  $\{\pi : \pi_k > 0, \sum_{k=1}^K \pi_k = 1\}$ .

- The posterior is also Dirichlet  $\text{Dir}(\alpha_1 + n_1, \dots, \alpha_K + n_K)$ .
- Posterior mean is

$$\hat{\pi}_k^{\text{mean}} = \frac{\alpha_k + n_k}{\sum_{j=1}^K \alpha_j + n_j}$$

## Dirichlet Distributions



- (A) Support of the Dirichlet density for  $K = 3$ .  
 (B) Dirichlet density for  $\alpha_k = 10$ .  
 (C) Dirichlet density for  $\alpha_k = 0.1$ .

## Naïve Bayes

- Return to the spam classification example with two-class naïve Bayes

$$p(x_i | \phi_k) = \prod_{j=1}^p \phi_{kj}^{x_i^{(j)}} (1 - \phi_{kj})^{1-x_i^{(j)}}.$$

- Set  $n_k = \sum_{i=1}^n \mathbb{1}\{y_i = k\}$ ,  $n_{kj} = \sum_{i=1}^n \mathbb{1}\{y_i = k, x_i^{(j)} = 1\}$ . MLE is:

$$\hat{\pi}_k = \frac{n_k}{n}, \quad \hat{\phi}_{kj} = \frac{\sum_{i: y_i=k} x_i^{(j)}}{n_k} = \frac{n_{kj}}{n_k}.$$

- One problem: if the  $\ell$ -th word did not appear in documents labelled as class  $k$  then  $\hat{\phi}_{k\ell} = 0$  and

$$\mathbb{P}(Y = k | X = x \text{ with } \ell\text{-th entry equal to } 1)$$

$$\propto \hat{\pi}_k \prod_{j=1}^p (\hat{\phi}_{kj})^{x^{(j)}} (1 - \hat{\phi}_{kj})^{1-x^{(j)}} = 0$$

i.e. we will never attribute a new document containing word  $\ell$  to class  $k$  (regardless of other words in it).

## Bayesian Inference on Naïve Bayes model

- Given  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ , want to predict a label  $\tilde{y}$  for a new document  $\tilde{x}$ . We can calculate

$$p(\tilde{x}, \tilde{y} = k | \mathcal{D}) = p(\tilde{y} = k | \mathcal{D}) p(\tilde{x} | \tilde{y} = k, \mathcal{D})$$

with

$$p(\tilde{y} = k | \mathcal{D}) = \frac{\alpha_k + n_k}{\sum_{l=1}^K \alpha_l + n}$$

$$p(\tilde{x}^{(j)} = 1 | \tilde{y} = k, \mathcal{D}) = \frac{a + n_{kj}}{a + b + n_k}$$

- Predicted class is

$$p(\tilde{y} = k | \tilde{x}, \mathcal{D}) = \frac{p(\tilde{y} = k | \mathcal{D}) p(\tilde{x} | \tilde{y} = k, \mathcal{D})}{p(\tilde{x} | \mathcal{D})}$$

- Compared to ML plug-in estimator, pseudocounts help to “regularize” probabilities away from extreme values.

## Bayesian Inference on Naïve Bayes model

- Under the Naïve Bayes model, the joint distribution of labels  $y_i \in \{1, \dots, K\}$  and data vectors  $x_i \in \{0, 1\}^p$  is

$$\prod_{i=1}^n p(x_i, y_i) = \prod_{i=1}^n \prod_{k=1}^K \left( \pi_k \prod_{j=1}^p \phi_{kj}^{x_i^{(j)}} (1 - \phi_{kj})^{1-x_i^{(j)}} \right)^{\mathbb{1}\{y_i=k\}}$$

$$= \prod_{k=1}^K \pi_k^{n_k} \prod_{j=1}^p \phi_{kj}^{n_{kj}} (1 - \phi_{kj})^{n_k - n_{kj}}$$

where  $n_k = \sum_{i=1}^n \mathbb{1}\{y_i = k\}$ ,  $n_{kj} = \sum_{i=1}^n \mathbb{1}\{y_i = k, x_i^{(j)} = 1\}$ .

- For conjugate prior, we can use  $\text{Dir}((\alpha_k)_{k=1}^K)$  for  $\pi$ , and  $\text{Beta}(a, b)$  for  $\phi_{kj}$  independently.
- Because the likelihood factorizes, the posterior distribution over  $\pi$  and  $(\phi_{kj})$  also factorizes, and posterior for  $\pi$  is  $\text{Dir}((\alpha_k + n_k)_{k=1}^K)$ , and for  $\phi_{kj}$  is  $\text{Beta}(a + n_{kj}, b + n_k - n_{kj})$ .

## Bayesian Learning and Regularization

- Consider a Bayesian approach to logistic regression: introduce a multivariate normal prior for weight vector  $w \in \mathbb{R}^p$ , and a uniform (improper) prior for offset  $b \in \mathbb{R}$ . The prior density is:

$$p(b, w) = 1 \cdot (2\pi\sigma^2)^{-\frac{p}{2}} \exp\left(-\frac{1}{2\sigma^2} \|w\|_2^2\right)$$

- The posterior is

$$p(b, w | \mathcal{D}) \propto \exp\left(-\frac{1}{2\sigma^2} \|w\|_2^2 - \sum_{i=1}^n \log(1 + \exp(-y_i(b + w^\top x_i)))\right)$$

- The posterior mode is equivalent to minimizing the  $L_2$ -regularized empirical risk.
- Regularized empirical risk minimization is (often) equivalent to having a prior and finding a MAP estimate of the parameters.
  - $L_2$  regularization - multivariate normal prior.
  - $L_1$  regularization - multivariate Laplace prior.
- From a Bayesian perspective, the MAP parameters are just one way to summarize the posterior distribution.

## Bayesian Model Selection

- A model  $\mathcal{M}$  with a given set of parameters  $\theta_{\mathcal{M}}$  consists of both the likelihood  $p(\mathcal{D}|\theta_{\mathcal{M}})$  and the prior distribution  $p(\theta_{\mathcal{M}})$ .
  - One example model would consist of all Gaussian mixtures with  $K$  components and equal covariance (LDA):  $\theta_{\text{LDA}} = (\pi_1, \dots, \pi_K; \mu_1, \dots, \mu_K; \Sigma)$ , along with a prior on  $\theta$ ; another would allow different covariances (QDA)  $\theta_{\text{QDA}} = (\pi_1, \dots, \pi_K; \mu_1, \dots, \mu_K; \Sigma_1, \dots, \Sigma_K)$ .
- The posterior distribution

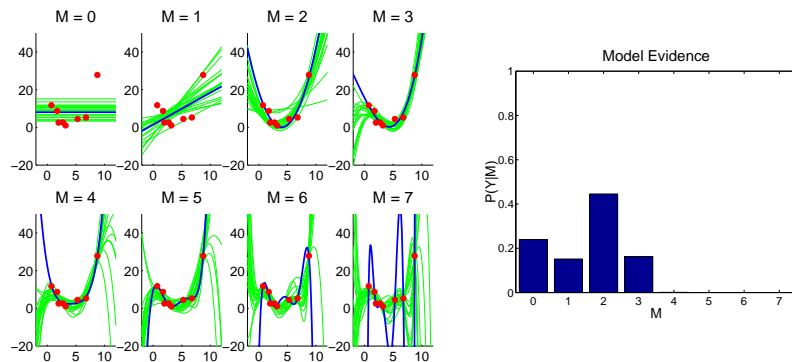
$$p(\theta_{\mathcal{M}}|\mathcal{D}, \mathcal{M}) = \frac{p(\mathcal{D}|\theta_{\mathcal{M}}, \mathcal{M})p(\theta_{\mathcal{M}}|\mathcal{M})}{p(\mathcal{D}|\mathcal{M})}$$

- Marginal probability of the data under  $\mathcal{M}$  (**Bayesian model evidence**):

$$p(\mathcal{D}|\mathcal{M}) = \int_{\Theta} p(\mathcal{D}|\theta_{\mathcal{M}}, \mathcal{M})p(\theta_{\mathcal{M}}|\mathcal{M})d\theta$$

- Compare models using their **Bayes factors**  $\frac{p(\mathcal{D}|\mathcal{M})}{p(\mathcal{D}|\mathcal{M}')}$

### Bayesian model comparison: Occam's razor at work



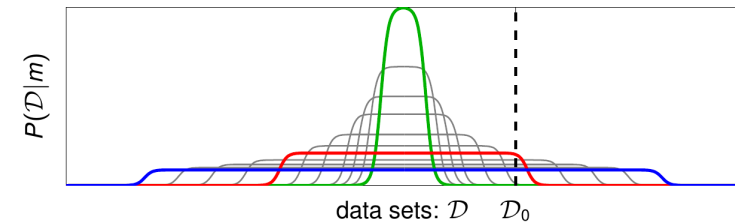
figures by M.Sahani

## Bayesian Occam's Razor

- **Occam's Razor**: of two explanations adequate to explain the same set of observations, the simpler should be preferred.

$$p(\mathcal{D}|\mathcal{M}) = \int_{\Theta} p(\mathcal{D}|\theta_{\mathcal{M}}, \mathcal{M})p(\theta_{\mathcal{M}}|\mathcal{M})d\theta$$

- Model evidence  $p(\mathcal{D}|\mathcal{M})$  is the probability that a set of randomly selected parameter values inside the model would generate dataset  $\mathcal{D}$ .
- Models that are **too simple** are unlikely to generate the observed dataset.
- Models that are **too complex** can generate many possible dataset, so again, they are unlikely to generate that particular dataset at random.



## Bayesian Learning – Discussion

- Use probability distributions to reason about uncertainties of parameters (latent variables and parameters are treated in the same way).
- Model consists of the likelihood function **and** the prior distribution on parameters: allows to integrate prior beliefs and domain knowledge.
- Bayesian computation — most posteriors are intractable, and posterior needs to be approximated by:
  - Monte Carlo methods (MCMC and SMC).
  - Variational methods (variational Bayes, belief propagation etc).
- Prior usually has hyperparameters, i.e.,  $p(\theta) = p(\theta|\psi)$ . How to choose  $\psi$ ?
  - Be Bayesian about  $\psi$  as well — choose a hyperprior  $p(\psi)$  and compute  $p(\psi|\mathcal{D})$ .
  - Maximum Likelihood II — find  $\psi$  maximizing  $\text{argmax}_{\psi \in \Psi} p(\mathcal{D}|\psi)$ .

$$p(\mathcal{D}|\psi) = \int p(\mathcal{D}|\theta)p(\theta|\psi)d\theta$$

$$p(\psi|\mathcal{D}) = \frac{p(\mathcal{D}|\psi)p(\psi)}{p(\mathcal{D})}$$

## Bayesian Learning – Further Reading

- Videlectures by Zoubin Ghahramani:  
Bayesian Learning and Graphical models.
- Gelman et al. Bayesian Data Analysis.
- Kevin Murphy. Machine Learning: a Probabilistic Perspective.
- E. T. Jaynes. Probability Theory: The Logic of Science.