

Bayesian Approach to Machine Learning (4)

Rohit Budhiraja

Machine Learning for Wireless Communications (EE798L)

Feb. 2, 2024

Recap of last lecture and today's agenda

- Recap of last class
 - Finished discussing Bayesian framework for coin toss example
- Today's agenda
 - Bayesian framework for Olympic data

Bayesian treatment for coin tossing game (recap)

- Modelled data using binomial distribution with likelihood

$$P(Y = y|r, N) = \binom{N}{y} r^y (1-r)^{N-y}$$

- Calculated winning probability with a ML (point) estimate of r i.e., $P(Y_{new} \leq 6|\hat{r}_{ML})$
- Considered r as a random variable, and captured its uncertainty
- By defining random variable Y_N as number of heads obtained in N tosses, we calculated $p(r|y_N)$

$$p(r|y_N) = \frac{P(y_N|r) p(r)}{P(y_N)}$$

- Re-calculated winning probability with MAP estimate of r i.e., $P(Y_{new} \leq 6|\hat{r}_{MAP})$
- Computed expected winning probability using **all of posterior information as well**

$$\int_{r=0}^{r=1} P(Y_{new} \leq 6|r) p(r|y_N) dr = \mathbb{E}_{p(r|y_N)}\{P(Y_{new} \leq 6|r)\}$$

- Computed marginal likelihood, and briefly informed that it captures model information

Maximum likelihood approach for Olympic data (recap)

- Recall our Olympic data model: $t_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$ such that

$$t_1 = w_0 + x_1 w_1 + \epsilon_1$$

$$\vdots = \vdots$$

$$t_N = w_0 + x_N w_1 + \epsilon_N$$

$$\mathbf{t} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

- Joint Gaussian likelihood of N data points

$$L = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\mathbf{w}^T \mathbf{x}_n, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N)$$

- Calculated **point ML estimate of \mathbf{w}** by maximizing logarithm likelihood which is $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$
- Predicted $t_{new} = \hat{\mathbf{w}}^T \mathbf{x}_{new}$, and calculated **predictive variance** $\sigma_{new}^2 = \text{var}\{t_{new}\}$
- Bayesian approach: calculate posterior distribution of \mathbf{w} , and use it to calculate t_{new}

Outline of Bayesian treatment of Olympic data

- We will use k th order polynomial model to model Olympic data

$$t_n = w_0 + w_1 x_n + w_2 x_n^2 + \dots + w_K x_n^K + \epsilon_n = \mathbf{w}^T \mathbf{x}_n + \epsilon_n$$

where $\mathbf{w} = [w_0, \dots, w_K]^T$ and $\mathbf{x}_n = [1, x_n, x_n^2, \dots, x_n^K]^T$, $\epsilon \sim \mathcal{N}(0, \sigma^2)$

- We can also write

$$\mathbf{t} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$$

where $\mathbf{t} = [t_1, \dots, t_N]^T$, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$ and $\boldsymbol{\epsilon} = [\epsilon_1, \dots, \epsilon_N]^T$

- Bayesian treatment enables us to incorporate insights about \mathbf{w} using a prior
 - Will define prior over \mathbf{w} using set of parameters Δ
- We assume that true value of σ^2 is known – simplify math
- Characterize randomness of \mathbf{w} by calculating posterior distribution using Bayes rule

$$p(\mathbf{w}|\mathbf{t}) = \frac{p(\mathbf{t}|\mathbf{w}) p(\mathbf{w})}{p(\mathbf{t})}$$

- Posterior of \mathbf{w} with dependence on parameter and data made explicit

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2, \Delta) p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)} = \frac{p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) p(\mathbf{w}|\Delta)}{p(\mathbf{t}|\mathbf{X}, \sigma^2, \Delta)}$$

Olympic data – choice of prior

- For our model $\mathbf{t} = \mathbf{X}\mathbf{w} + \epsilon$, likelihood is Gaussian $p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) = \mathcal{N}(\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N)$
- Use a Gaussian prior $p(\mathbf{w}|\Delta)$, which is conjugate to Gaussian likelihood,

$$p(\mathbf{w}|\mu_0, \Sigma_0) = \mathcal{N}(\mu_0, \Sigma_0),$$

- We will choose parameters μ_0, Σ_0 later. Also, we will not always explicitly condition on μ_0, Σ_0
- For example, instead of $p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \Delta) = p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2, \mu_0, \Sigma_0)$ we will use $p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2)$

Olympic data – Posterior calculation (1)

- We will use the fact that posterior will be Gaussian – allows us to ignore the marginal likelihood

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) &\propto p(\mathbf{t}|\mathbf{w}, \mathbf{X}, \sigma^2) p(\mathbf{w}|\mu_0, \Sigma_0) \\ &= \frac{1}{(2\pi)^{N/2} |\sigma^2 \mathbf{I}|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\sigma^2 \mathbf{I})^{-1} (\mathbf{t} - \mathbf{X}\mathbf{w})\right) \\ &\times \frac{1}{(2\pi)^{N/2} |\Sigma_0|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{w} - \mu_0)^T \Sigma_0^{-1} (\mathbf{w} - \mu_0)\right) \\ &\stackrel{(a)}{\propto} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})\right) \times \exp\left(-\frac{1}{2} (\mathbf{w} - \mu_0)^T \Sigma_0^{-1} (\mathbf{w} - \mu_0)\right) \\ &= \exp\left\{-\frac{1}{2} \left(-\frac{1}{\sigma^2} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w}) + (\mathbf{w} - \mu_0)^T \Sigma_0^{-1} (\mathbf{w} - \mu_0)\right)\right\}. \end{aligned}$$

- Proportionality (a) is because we ignore the terms which are independent of \mathbf{w}
- Multiplying the terms in bracket, and once again removing any that don't involve \mathbf{w} gives

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) \propto \exp\left\{-\frac{1}{2} \left(-\frac{2}{\sigma^2} \mathbf{t}^T \mathbf{X}\mathbf{w} + \frac{1}{\sigma^2} \mathbf{w}^T \mathbf{X}^T \mathbf{X}\mathbf{w} + \mathbf{w}^T \Sigma_0^{-1} \mathbf{w} - 2\mu_0^T \Sigma_0^{-1} \mathbf{w}\right)\right\}$$

Olympic data – Posterior calculation (2)

- We have

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) \propto \exp \left\{ -\frac{1}{2} \left(-\frac{2}{\sigma^2} \mathbf{t}^T \mathbf{X} \mathbf{w} + \frac{1}{\sigma^2} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \boldsymbol{\Sigma}_0^{-1} \mathbf{w} - 2\boldsymbol{\mu}_0^T \boldsymbol{\Sigma}_0^{-1} \mathbf{w} \right) \right\} \quad (1)$$

- We take a generic multivariate Gaussian pdf and rearrange it to make it look like (1)

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) &= \mathcal{N}(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \\ &\propto \exp \left(-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu}_w)^T \boldsymbol{\Sigma}_w^{-1} (\mathbf{w} - \boldsymbol{\mu}_w) \right) \\ &\propto \exp \left\{ -\frac{1}{2} (\mathbf{w}^T \boldsymbol{\Sigma}_w^{-1} \mathbf{w} - 2\boldsymbol{\mu}_w^T \boldsymbol{\Sigma}_w^{-1} \mathbf{w}) \right\} \end{aligned} \quad (2)$$

- Linear and quadratic term in \mathbf{w} in (1) must be equal to those in (2). Start with quadratic

$$\begin{aligned} \mathbf{w}^T \boldsymbol{\Sigma}_w^{-1} \mathbf{w} &= \frac{1}{\sigma^2} \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \boldsymbol{\Sigma}_0^{-1} \mathbf{w} = \mathbf{w}^T \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right) \mathbf{w} \\ \boldsymbol{\Sigma}_w &= \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \end{aligned}$$

Olympic data – Posterior calculation (3)

- Similarly, equating linear terms from posterior we can get an expression for μ_w :

$$\begin{aligned}-2\mu_w^T \Sigma_w^{-1} w &= -\frac{2}{\sigma^2} \mathbf{t}^T \mathbf{X} w - 2\mu_0^T \Sigma_0^{-1} w \\ \mu_w^T \Sigma_w^{-1} w &= \frac{1}{\sigma^2} \mathbf{t}^T \mathbf{X} w + \mu_0^T \Sigma_0^{-1} w \\ \mu_w^T \Sigma_w^{-1} &= \frac{1}{\sigma^2} \mathbf{t}^T \mathbf{X} + \mu_0^T \Sigma_0^{-1} \\ \mu_w^T \Sigma_w^{-1} \Sigma_w &= \left(\frac{1}{\sigma^2} \mathbf{t}^T \mathbf{X} + \mu_0^T \Sigma_0^{-1} \right) \Sigma_w \\ \mu_w^T &= \left(\frac{1}{\sigma^2} \mathbf{t}^T \mathbf{X} + \mu_0^T \Sigma_0^{-1} \right) \Sigma_w\end{aligned}$$

- Using the fact $\Sigma_w^T = \Sigma_w$, we have

$$\mu_w = \Sigma_w \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t} + \Sigma_0^{-1} \mu_0 \right),$$

Olympic data – Posterior calculation (4)

- Posterior is therefore

$$p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}})$$

with

$$\boldsymbol{\Sigma}_{\mathbf{w}} = \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1}$$

and

$$\boldsymbol{\mu}_{\mathbf{w}} = \boldsymbol{\Sigma}_{\mathbf{w}} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t} + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right)$$

- Mean (mode) $\hat{\mathbf{w}} = \boldsymbol{\mu}_{\mathbf{w}}$ is (MAP) estimate. If prior has zero mean $\boldsymbol{\mu}_0 = \mathbf{0}$ then MAP estimate is

$$\boldsymbol{\mu}_{\mathbf{w}} = \boldsymbol{\Sigma}_{\mathbf{w}} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{t} \right) = \frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} \mathbf{X}^T \mathbf{X} + \boldsymbol{\Sigma}_0^{-1} \right)^{-1} \mathbf{X}^T \mathbf{t}$$

- Similar to regularized least squares estimate
- Given a new observation \mathbf{x}_{new} , predict t_{new} using point estimate
- Treat t_{new} as a random variable by casting $t_{new} = \mathbf{w}^T \mathbf{x}_{new} + \epsilon_{new}$

$$p(t_{new}|\mathbf{x}_{new}, \hat{\mathbf{w}}, \sigma^2) = \mathcal{N}(\hat{\mathbf{w}}^T \mathbf{x}_{new}, \sigma^2)$$

- Using MAP estimate, $t_{new} = \hat{\mathbf{w}}^T \mathbf{x}_{new}$, and has variance σ^2

Predictive distribution calculation (summary)

- Predict t_{new} by capturing complete randomness in \mathbf{w} using its posterior distribution

$$\int p(t_{new}|\mathbf{x}_{new}, \mathbf{w}, \sigma^2) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) d\mathbf{w} = p(t_{new}|\mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2)$$

- $p(t_{new}|\cdot)$, also called **predictive distribution**, and is not conditioned on \mathbf{w}
- Recall $t_{new} = \mathbf{w}^T \mathbf{x}_{new} + \epsilon_{new}$ with

$$p(t_{new}|\mathbf{x}_{new}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{x}_{new}^T \mathbf{w}, \sigma^2) \text{ with } p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) = \mathcal{N}(\boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}})$$

Theorem

If marginal and conditional distributions of a generic Gaussian vector \mathbf{x} are

$$p(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \mathbf{\Lambda}^{-1}) \text{ and } p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}).$$

Here $\mathbf{\Lambda}$ and \mathbf{L} are precision (inverse of covariance) matrices. Marginal distribution of \mathbf{y} is

$$p(\mathbf{y}) = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x})d\mathbf{x} = \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\mathbf{\Lambda}^{-1}\mathbf{A}^T)$$

- Assume $\mathbf{w} = \mathbf{x}$ and $t_{new} = y$ and with $\boldsymbol{\mu} = \boldsymbol{\mu}_{\mathbf{w}}$, $\mathbf{\Lambda}^{-1} = \boldsymbol{\Sigma}_{\mathbf{w}}$, $b = 0$, $\mathbf{L}^{-1} = \sigma^2$, $\mathbf{A} = \mathbf{x}_{new}^T$, we have

$$p(t_{new}|\mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2) = \int p(t_{new}|\mathbf{x}_{new}, \mathbf{w}, \sigma^2) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \sigma^2) d\mathbf{w} = \mathcal{N}(\mathbf{x}_{new}^T \boldsymbol{\mu}_{\mathbf{w}}, \sigma^2 + \mathbf{x}_{new}^T \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x}_{new})$$

- $t_{new} = \mathbf{x}_{new}^T \boldsymbol{\mu}_{\mathbf{w}}$ with variance $\sigma^2 + \mathbf{x}_{new}^T \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x}_{new}$

Predicting using MAP estimate of \mathbf{w} and its complete posterior

- Given a new observation \mathbf{x}_{new} , we want to predict t_{new}
- Treat t_{new} as a random variable by casting $t_{new} = \mathbf{w}^T \mathbf{x}_{new} + \epsilon_{new}$

$$p(t_{new} | \mathbf{x}_{new}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{w}^T \mathbf{x}_{new}, \sigma^2) \quad (3)$$

- With MAP estimate $\hat{\mathbf{w}}$, Eq. (3) using $\hat{\mathbf{w}}$, is $p(t_{new} | \mathbf{x}_{new}, \hat{\mathbf{w}}, \sigma^2) = \mathcal{N}(\hat{\mathbf{w}}^T \mathbf{x}_{new}, \sigma^2)$
 - $t_{new} = \hat{\mathbf{w}}^T \mathbf{x}_{new} = \boldsymbol{\mu}_{\mathbf{w}}^T \mathbf{x}_{new}$ with variance σ^2
- With complete distribution of \mathbf{w} , we calculated predictive distribution

$$\int p(t_{new} | \mathbf{x}_{new}, \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathbf{t}, \mathbf{X}, \sigma^2) d\mathbf{w} = p(t_{new} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2)$$

- We showed $p(t_{new} | \mathbf{x}_{new}, \mathbf{X}, \mathbf{t}, \sigma^2) = \mathcal{N}(\mathbf{x}_{new}^T \boldsymbol{\mu}_{\mathbf{w}}, \sigma^2 + \mathbf{x}_{new}^T \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x}_{new})$
 - Implies $t_{new} = \mathbf{x}_{new}^T \boldsymbol{\mu}_{\mathbf{w}}$ with variance $\sigma^2 + \mathbf{x}_{new}^T \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x}_{new}$