

Bayesian Linear Regression (Contd)

Piyush Rai

Topics in Probabilistic Modeling and Inference (CS698X)

Jan 16, 2018

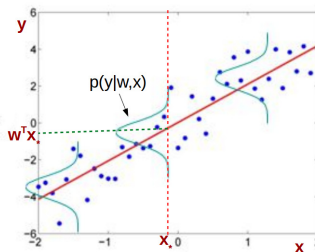
Linear Regression: A Probabilistic Setup

- Given: N training examples $\{\mathbf{x}_n, y_n\}_{n=1}^N$, features: $\mathbf{x}_n \in \mathbb{R}^D$, response $y_n \in \mathbb{R}$
- Assume a “noisy” linear model with regression weight vector $\mathbf{w} = [w_1, w_2, \dots, w_D] \in \mathbb{R}^D$

$$y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$$

where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$, β : precision (inverse variance) of Gaussian (assumed known)

- Therefore $p(y_n | \mathbf{x}_n, \mathbf{w}, \beta) = \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$



- Note: Some books (e.g., PRML) use $\phi(\mathbf{x}_n)$ to denote the features where ϕ is some transformation of the original features \mathbf{x}_n (we will only use this notation when talking about nonlinear regression)

The Likelihood Model

- Notation: $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$: $N \times D$ feature matrix, $\mathbf{y} = [y_1 \dots y_N]^\top$: $N \times 1$ response vector
- Assuming i.i.d. observations, the likelihood model

$$\begin{aligned} p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta) &= \prod_{n=1}^N p(y_n|\mathbf{w}, \mathbf{x}_n, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|\mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) \\ &= \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} \exp\left[-\frac{\beta}{2}(y_n - \mathbf{w}^\top \mathbf{x}_n)^2\right] \\ &= \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \exp\left[-\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2\right] \end{aligned}$$

- Can also write the likelihood $p(\mathbf{y}|\mathbf{w}, \mathbf{X})$ as an N -dim multivariate Gaussian

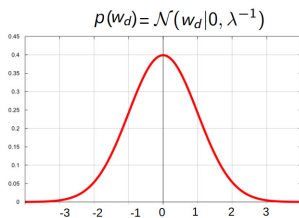
$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \exp\left[-\frac{\beta}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\right]$$

- Note that NLL = sum of squared errors! Minimizing w.r.t. \mathbf{w} will give MLE/least squares solution!

The Prior

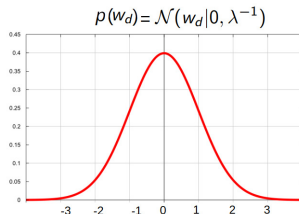
- Assume the entries in \mathbf{w} are i.i.d. with zero mean Gaussian priors. Therefore

$$p(\mathbf{w}) = \prod_{d=1}^D p(w_d) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \left(\frac{\lambda}{2\pi}\right)^{\frac{D}{2}} \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right]$$



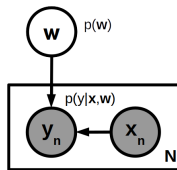
- This prior promotes the entries in \mathbf{w} to be small (close to zero)
 - Also, the negative of log-prior is the same as an ℓ_2 regularizer on \mathbf{w}
- This prior is conjugate to the likelihood (Gaussian) which makes posterior inference easy

The Prior



- The role of the precision hyperparam λ in the prior is important
- Large values of λ would more aggressively encourage w_d to be close to zero
- Can think of λ as the regularization hyperparam for the weights
- **Important:** Can infer λ as well (will see later how to do this)
- Can even have different λ for each w_d , i.e., $p(\mathbf{w} | \{\lambda_d\}_{d=1}^D) = \prod_{d=1}^D \mathcal{N}(w_d | 0, \lambda_d^{-1})$
 - Useful in **sparse regression/classification** models in which very few features are relevant which can be identified by inferring $\{\lambda_d\}_{d=1}^D$. More on this when we talk about sparse models.

Inference Tasks for Bayesian Linear Regression



(Hyperparameters λ, β not shown as they are fixed/known)

- Want to infer the posterior distribution over \mathbf{w} (for now, assume β and λ to be known)

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}, \beta, \lambda) = \frac{p(\mathbf{w} | \lambda) p(\mathbf{y} | \mathbf{w}, \mathbf{X}, \beta)}{p(\mathbf{y} | \mathbf{X}, \beta, \lambda)}$$

- Want to infer the posterior predictive distribution

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \int p(y_* | \mathbf{w}, \mathbf{x}_*, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) d\mathbf{w}$$

- Since the likelihood model $p(\mathbf{y} | \mathbf{w}, \mathbf{x}, \beta)$ and the prior $p(\mathbf{w} | \lambda)$ both are Gaussians, the above computations can be easily done

Brief Detour: Some Useful Properties of Multivariate Gaussians

Multivariate Gaussian Distribution

- The (multivariate) Gaussian with mean μ and cov. matrix Σ

$$\begin{aligned}\mathcal{N}(\mathbf{x}|\mu, \Sigma) &= \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu)^\top \Sigma^{-1} (\mathbf{x} - \mu) \right\} \\ &= \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} \exp \left\{ -\frac{1}{2} \text{trace} \left[\Sigma^{-1} \mathbf{S} \right] \right\} \quad \text{where } \mathbf{S} = (\mathbf{x} - \mu)(\mathbf{x} - \mu)^\top\end{aligned}$$

- An alternate representation: The “information form”

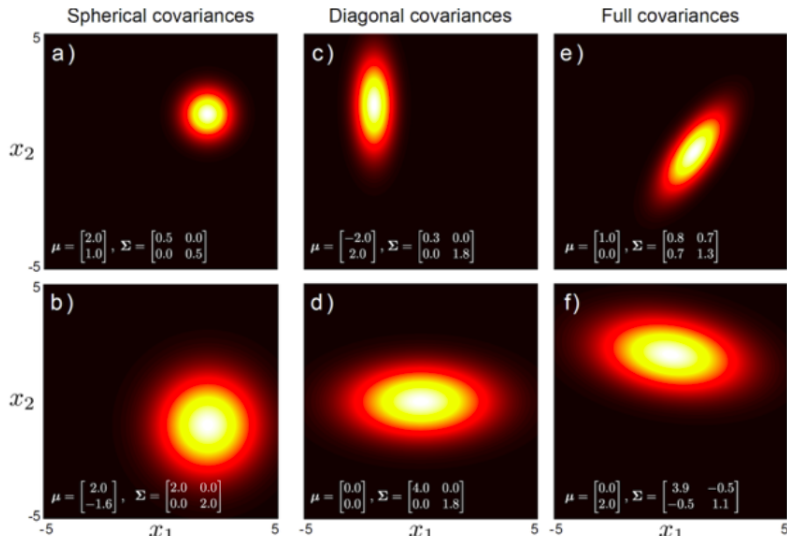
$$\mathcal{N}_c(\mathbf{x}|\xi, \Lambda) = (2\pi)^{-D/2} |\Lambda|^{1/2} \exp \left\{ -\frac{1}{2} \left(\mathbf{x}^\top \Lambda \mathbf{x} + \xi^\top \Lambda^{-1} \xi - 2\mathbf{x}^\top \xi \right) \right\}$$

where $\Lambda = \Sigma^{-1}$ (precision matrix) and $\xi = \Sigma^{-1} \mu$ are the “natural parameters”

- Note that there is a term **quadratic in \mathbf{x}** (involves $\Lambda = \Sigma^{-1}$) and **linear in \mathbf{x}** (involves $\xi = \Sigma^{-1} \mu$)
- Information form can help recognize μ and Σ of a Gaussian when doing algebraic manipulations

Multivariate Gaussian: The Covariance Matrix

The covariance matrix can be spherical, diagonal, or full



Marginals and Conditionals from Gaussian Joint Distribution

- Given \mathbf{x} having multivariate Gaussian distribution $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$. Suppose

$$\begin{aligned}\mathbf{x} &= \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix} & \boldsymbol{\mu} &= \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix} \\ \boldsymbol{\Sigma} &= \begin{bmatrix} \boldsymbol{\Sigma}_{aa} & \boldsymbol{\Sigma}_{ab} \\ \boldsymbol{\Sigma}_{ba} & \boldsymbol{\Sigma}_{bb} \end{bmatrix} & \boldsymbol{\Lambda} &= \begin{bmatrix} \boldsymbol{\Lambda}_{aa} & \boldsymbol{\Lambda}_{ab} \\ \boldsymbol{\Lambda}_{ba} & \boldsymbol{\Lambda}_{bb} \end{bmatrix}\end{aligned}$$

- The **marginal distribution** of one block, say \mathbf{x}_a , is a Gaussian

$$p(\mathbf{x}_a) = \int p(\mathbf{x}_a, \mathbf{x}_b) d\mathbf{x}_b = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

- The **conditional distribution** of \mathbf{x}_a given \mathbf{x}_b , is Gaussian, i.e., $p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, \boldsymbol{\Sigma}_{a|b})$ where

$$\begin{aligned}\boldsymbol{\Sigma}_{a|b} &= \boldsymbol{\Lambda}_{aa}^{-1} = \boldsymbol{\Sigma}_{aa} - \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} \boldsymbol{\Sigma}_{ba} & (\text{"smaller" than } \boldsymbol{\Sigma}_{aa}; \text{ makes sense intuitively}) \\ \boldsymbol{\mu}_{a|b} &= \boldsymbol{\Sigma}_{a|b} \{ \boldsymbol{\Lambda}_{aa} \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \} \\ &= \boldsymbol{\mu}_a - \boldsymbol{\Lambda}_{aa}^{-1} \boldsymbol{\Lambda}_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_{ab} \boldsymbol{\Sigma}_{bb}^{-1} (\mathbf{x}_b - \boldsymbol{\mu}_b)\end{aligned}$$

- Both results are **extremely useful** especially when working with Gaussian joint distributions

Linear Gaussian Model

- Consider **linear transformation** of a Gaussian r.v. \mathbf{z} with $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$, plus **Gaussian noise**

$$\boxed{\mathbf{x} = \mathbf{A}\mathbf{z} + \mathbf{b} + \boldsymbol{\epsilon}} \quad \text{where} \quad p(\boldsymbol{\epsilon}) = \mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{L}^{-1})$$

- Easy to see that, conditioned on \mathbf{z} , \mathbf{x} too has a Gaussian distribution

$$p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\mathbf{A}\mathbf{z} + \mathbf{b}, \mathbf{L}^{-1})$$

- This is called a **Linear Gaussian Model**. Very commonly encountered in probabilistic modeling
- The following two distributions are of particular interest. Defining $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$, we have

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})} = \mathcal{N}(\mathbf{z}|\boldsymbol{\Sigma} \{ \mathbf{A}^\top \mathbf{L}(\mathbf{x} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu} \}, \boldsymbol{\Sigma}) \quad (\text{Gaussian posterior})$$

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{x}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^\top + \mathbf{L}^{-1}) \quad (\text{Gaussian predictive/marginal})$$

- Exercise:** Prove the above two results (MLAPP Chap. 4 and PRML Chap. 2 contain the proof)
 - Can derive by first writing down the joint $p(\mathbf{x}, \mathbf{z})$, identifying its mean and covariance/precision matrix (using information form), and then apply conditioning and marginal formulas from previous slide

Bayesian Linear Regression: The Posterior

- The posterior over \mathbf{w} (for now, assume hyperparams β and λ to be known)

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \frac{p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta)}{p(\mathbf{y}|\mathbf{X}, \beta, \lambda)} \propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \beta)$$

- Computing $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) \propto \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) \times \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N)$$

- Using the “completing the squares” trick (or directly using Gaussian conditioning formula)

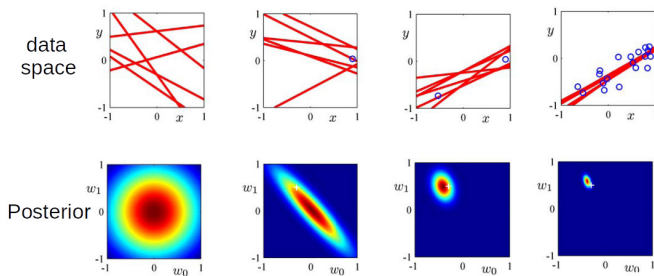
$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

$$\text{where } \boldsymbol{\Sigma}_N = \left(\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D \right)^{-1} = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \quad (\text{posterior's covariance matrix})$$

$$\boldsymbol{\mu}_N = \boldsymbol{\Sigma}_N \left[\beta \sum_{n=1}^N y_n \mathbf{x}_n \right] = \boldsymbol{\Sigma}_N [\beta \mathbf{X}^\top \mathbf{y}] = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y} \quad (\text{posterior's mean})$$

The Posterior: A Visualization

- Assume a linear regression problem with ground truth $\mathbf{w} = [w_0, w_1]$ with $w_0 = -0.3$, $w_1 = 0.5$
- Assume data generated by a linear regression model $y = w_0 + w_1x + \text{"noise"}$
 - Note: It's actually 1-D regression (w_0 is just a bias term), or 2-D reg. with feature $[1, x]$
- Figures below show the “data space” and posterior of \mathbf{w} for different number of observations (note: with no observations, the posterior = prior)



- The “data space” (red lines) shown above denotes various possible linear regression datasets with data of the form $y = w_0 + w_1x$ generated using \mathbf{w} drawn from the current posterior of \mathbf{w}

Bayesian Linear Regression: Posterior Predictive Distribution

- Given the posterior $p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$, how to make prediction y_* for a new input \mathbf{x}_* ?
- The posterior predictive distribution will be

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \int p(y_*|\mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) d\mathbf{w}$$

- Using Gaussian predictive/marginal formula, the posterior predictive will be another Gaussian

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*)$$

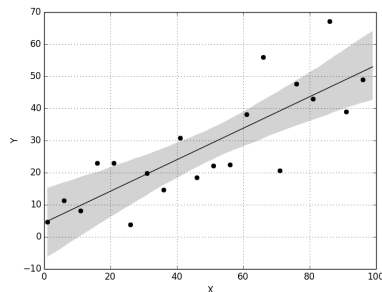
- So we get a predictive mean $\boldsymbol{\mu}_N^\top \mathbf{x}_*$ and an input-specific predictive variance $\beta^{-1} + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*$
- In contrast, MLE and MAP make “plug-in” predictions (using the point estimate of \mathbf{w})

$$\begin{aligned} p(y_*|\mathbf{x}_*, \mathbf{w}_{MLE}) &= \mathcal{N}(\mathbf{w}_{MLE}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MLE prediction} \\ p(y_*|\mathbf{x}_*, \mathbf{w}_{MAP}) &= \mathcal{N}(\mathbf{w}_{MAP}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MAP prediction} \end{aligned}$$

- Important: Unlike MLE/MAP, the variance of y_* also depends on the input \mathbf{x}_* (this, as we will see later, will be very useful in sequential decision-making problems such as active learning)

Posterior Predictive Distribution: An Illustration

Black dots are training examples



Width of the shaded region at any x denotes the predictive uncertainty at that x (\pm one std-dev)

Regions with more training examples have smaller predictive variance

Interpreting the Predictions

- We already saw that the posterior predictive for a new test point \mathbf{x}_* is a Gaussian

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}_N^\top \mathbf{x}_*, \beta^{-1} + \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_*)$$

- (Exercise) Can show that the mean prediction at a new test input \mathbf{x}_* can also be written as

$$\mathbb{E}_{p(\mathbf{w}|\mathbf{y}, \mathbf{X})}[\mathbf{w}^\top \mathbf{x}_*] = \boldsymbol{\mu}_N^\top \mathbf{x}_* = \sum_{n=1}^N k(\mathbf{x}_n, \mathbf{x}_*) y_n \quad (\text{where } k(\mathbf{x}_n, \mathbf{x}_*) = \beta \mathbf{x}_*^\top \boldsymbol{\Sigma}_N \mathbf{x}_n)$$

.. which makes intuitive sense

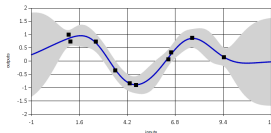
- (Exercise) Covariance of the predictions at any two inputs \mathbf{x} and \mathbf{x}' will be

$$\text{cov}[\mathbf{x}^\top \mathbf{w}, \mathbf{w}^\top \mathbf{x}'] = \mathbf{x}^\top \boldsymbol{\Sigma}_N \mathbf{x}' = \beta^{-1} k(\mathbf{x}, \mathbf{x}')$$

.. which makes intuitive sense (predictions of two faraway points will have low covariance)

- These interpretations will be useful when talking about Gaussian Processes for nonlinear regression

Nonlinear Regression?



- Can extend the linear regression model to handle nonlinear regression problems
- One way is to replace the feature vectors \mathbf{x} by a nonlinear mapping $\phi(\mathbf{x})$

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \phi(\mathbf{x}), \beta^{-1})$$

- The nonlinear mapping can be defined directly, e.g., for a one-dimensional feature x

$$\phi(x) = [1, x, x^2]$$

- Alternatively, a kernel function can be used to implicitly define the nonlinear mapping
- More on nonlinear regression when we discuss Gaussian Processes

What about the hyperparameters of the regression model?

- If hyperparameters are to be estimated, we will have a hierarchical/multiparameter model
- Posterior inference is slightly more involved in this case
- Iterative methods required to learn the weight vector and the hyperparameters, e.g.,
 - Marginal likelihood maximization for hyperparameter estimation
 - Expectation maximization (EM)
 - MCMC or variational inference
- We will discuss more when we talk about inference in hierarchical/multiparameter models