

Probabilistic Linear Regression

Piyush Rai
IIT Kanpur

Probabilistic Machine Learning (CS772A)

Jan 13, 2016

Linear Regression: A Probabilistic View

- Given: N training examples $\{\mathbf{x}_n, y_n\}_{n=1}^N$, features: $\mathbf{x}_n \in \mathbb{R}^D$, response $y_n \in \mathbb{R}$
- $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$: $N \times D$ feat. matrix, $\mathbf{Y} = [y_1 \dots y_N]^\top$: $N \times 1$ resp. vector
- Probabilistic view: responses are generated via a probabilistic model
- Assume a “noisy” linear model with regression weight vector $\mathbf{w} \in \mathbb{R}^D$:

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

- Gaussian noise**: $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$, β : precision (inverse variance) of Gaussian
- Thus each response y_n also has a Gaussian distribution

$$y_n \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$$

- Goal: Learn regression weight vector \mathbf{w} to predict y_* for a new \mathbf{x}_*

Linear Regression: A Probabilistic View

- For Gaussian response y_n

$$p(y_n | \mathbf{x}_n, \mathbf{w}) = \sqrt{\frac{\beta}{2\pi}} \exp \left\{ -\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right\}$$

- Thus the likelihood (assuming i.i.d. responses) or *probability* of data:

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \left(\frac{\beta}{2\pi} \right)^{\frac{N}{2}} \exp \left\{ -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right\}$$

- Note: \mathbf{x}_n (features) assumed given/fixed. Only modeling the response y_n
- Log-likelihood** (ignoring constants w.r.t. \mathbf{w})

$$\log p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) \propto -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Note that the log-likelihood is nothing but a (weighted) sum of (negative) squared errors on training data: high log-lik \Rightarrow low sum of squared errors

Maximum Likelihood Estimation (MLE)

- MLE: Find the \mathbf{w} that maximizes the (log) likelihood $\log p(\mathbf{Y} | \mathbf{X}, \mathbf{w})$

$$\arg \max_{\mathbf{w}} \log p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \arg \min_{\mathbf{w}} -\log p(\mathbf{Y} | \mathbf{X}, \mathbf{w}) = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Same objective as the classic **ordinary least squares** (OLS) regression
 - Basically, maximizing log-lik = minimizing the sum of squared errors
- Taking derivative w.r.t. \mathbf{w} and setting to zero, we get

$$\mathbf{w}_{MLE} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- Same solution as the solution of the OLS regression problem. Some issues:
 - $\mathbf{X}^\top \mathbf{X}$ may be ill-conditioned (not invertible)
 - “Uncontrolled” \mathbf{w} can lead to overfitting (thus need regularization)
- A solution: Put a **prior distribution** on \mathbf{w} (to impose “smoothness” and control \mathbf{w}) and do MAP estimation (MAP estimation = “regularized” MLE)

Prior Distribution on Weights

- Assume zero-mean spherical **Gaussian prior** on weights $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_D]$

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

λ is **precision** (inverse variance) of the Gaussian and $\|\mathbf{w}\|^2 = \sum_{d=1}^D w_d^2$

- Note: We can also write the prior as a product of D univariate Gaussians

$$p(\mathbf{w}) = \prod_{d=1}^D p(w_d) = \prod_{d=1}^D \mathcal{N}(0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} \exp\left(-\frac{\lambda}{2} w_d^2\right) = \left(\frac{\lambda}{2\pi}\right)^{D/2} \exp\left(-\frac{\lambda}{2} \sum_{d=1}^D w_d^2\right)$$

- Gaussian prior encourages a “small” \mathbf{w} by shrinking each component w_d towards zero (Gaussian’s mean). Precision λ controls the extent of shrinkage
- This corresponds to imposing a **regularizer** on \mathbf{w} . We will soon see (or you might already have guessed) that the **Gaussian prior** results in a **squared norm (ℓ_2) regularizer**, and λ controls the strength of regularization
- Note: Different types of priors result in different types of regularizers (e.g., a **Laplace prior** on \mathbf{w} : $p(\mathbf{w}) \propto \exp(-|\mathbf{w}|)$ will result in an ℓ_1 regularizer on \mathbf{w})

MAP Estimation

- The posterior distribution on \mathbf{w} : $p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \propto p(\mathbf{Y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})$
- The (log) posterior: $\log p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) = \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) + \log p(\mathbf{w})$. Thus,

$$\log p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \propto -\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 - \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w} \quad (\text{ignoring constants w.r.t } \mathbf{w})$$

- MAP Estimation: Maximize the (log) posterior w.r.t. \mathbf{w}

$$\arg \max_{\mathbf{w}} \log p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{w}} -\log p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) = \arg \min_{\mathbf{w}} \underbrace{\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2}_{\text{fit to the training data}} + \underbrace{\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}_{\text{keep } \mathbf{w} \text{ "simple"}}$$

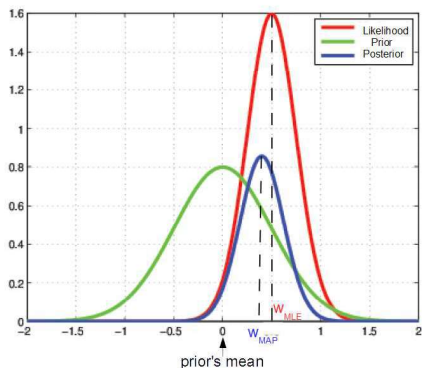
- Thus MAP estimation finds a \mathbf{w} by trying to balance between the **likelihood** (fit to the training data) vs the **prior** (model’s simplicity)
- Setting derivative w.r.t. \mathbf{w} to zero yields

$$\mathbf{w}_{MAP} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \frac{\lambda}{\beta} \mathbf{I}_D\right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{Y}$$

- This corresponds to the solution of the **ridge regression** (regularized least squares) problem with regularization parameter $\frac{\lambda}{\beta}$

MAP Estimation: An Illustration

\mathbf{w}_{MAP} is a compromise between prior’s mean and \mathbf{w}_{MLE}



Summary: MLE vs MAP for Linear Regression

- MLE Objective

$$\arg \max_{\mathbf{w}} \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}) = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- MLE solution

$$\mathbf{w}_{MLE} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top\right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$$

- MAP Objective

$$\arg \max_{\mathbf{w}} \log p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) \propto \arg \max_{\mathbf{w}} \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}) = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{\beta} \mathbf{w}^\top \mathbf{w}$$

- MAP solution

$$\mathbf{w}_{MAP} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \frac{\lambda}{\beta} \mathbf{I}_D\right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{Y}$$

The “Fully” Bayesian Approach

- MLE/MAP only provide a point estimate of \mathbf{w} (no estimate of uncertainty)
- Let's try to infer the full posterior of \mathbf{w} : $p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{p(\mathbf{Y}|\mathbf{X})}$
- Since the likelihood and the prior, both, are Gaussian, the posterior will also be Gaussian (due to conjugacy)
- What will be the posterior's mean and covariance/precision matrix ?
- Since \mathbf{X} is known/fixed, and using the property of Gaussians, given $p(\mathbf{Y}|\mathbf{X}, \mathbf{w})$ and $p(\mathbf{w})$ both Gaussian (refer to the results discussed in lecture 2),

$$p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\text{where } \boldsymbol{\mu} = \boldsymbol{\Sigma}(\beta \sum_{n=1}^N y_n \mathbf{x}_n) = \boldsymbol{\Sigma}(\beta \mathbf{X}^T \mathbf{Y})$$

$$\boldsymbol{\Sigma} = (\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T + \lambda \mathbf{I}_D)^{-1} = (\beta \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_D)^{-1}$$

Making Predictions

- MLE and MAP make “plug-in” predictions

$$p(y_*|\mathbf{x}_*, \mathbf{w}_{MLE}) = \mathcal{N}(\mathbf{w}_{MLE}^T \mathbf{x}_*, \beta^{-1}) \quad \text{-- MLE prediction}$$

$$p(y_*|\mathbf{x}_*, \mathbf{w}_{MAP}) = \mathcal{N}(\mathbf{w}_{MAP}^T \mathbf{x}_*, \beta^{-1}) \quad \text{-- MAP prediction}$$

- MLE/MAP only use a point estimate ($\mathbf{w}_{MLE}/\mathbf{w}_{MAP}$) for making prediction
- Fully Bayesian approach of making predictions is via the **predictive posterior**

$$p(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{Y}) = \int_{\mathbf{w}} p(y_*|\mathbf{x}_*, \mathbf{w}) p(\mathbf{w}|\mathbf{X}, \mathbf{Y}) d\mathbf{w} \quad (\text{Predictive Posterior})$$

- **Predictive Posterior:** Don't use a single \mathbf{w} to make predictions but average $p(y_*|\mathbf{x}_*, \mathbf{w})$ over all possible \mathbf{w} 's (each weighted by its posterior probability)
- Since the likelihood $p(y_*|\mathbf{x}_*, \mathbf{w})$ and posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$ are Gaussian, the predictive posterior is also Gaussian. Thus in the fully Bayesian approach:

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

where $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ are mean and cov. matrix, resp., of the posterior $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$

Some things we didn't cover..

- How to estimate the model hyperparameters (e.g., precisions β and λ)? The Bayesian approach allows us doing this.
- Nonlinear regression. What to do when a linear model doesn't fit the responses well. Kernel methods (e.g., Gaussian Processes) can handle this.

(We will see these later in the semester)

Next class: Logistic Regression