

Probabilistic Linear Regression

CS698X: Topics in Probabilistic Modeling and Inference

Piyush Rai

Probabilistic Linear Regression



Note: Only y_n being modeled, not x_n (discriminative model). A **conditional model** where y_n being modeled, conditioned on x_n

2

- Assume training data $\{\mathbf{x}_n, y_n\}_{n=1}^N$, with features $\mathbf{x}_n \in \mathbb{R}^D$ and responses $y_n \in \mathbb{R}$
- Assume each y_n generated by a noisy **linear model** with wts $\mathbf{w} = [w_1, \dots, w_D]$

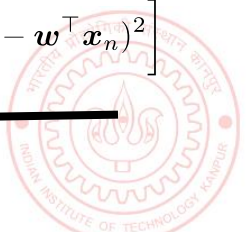
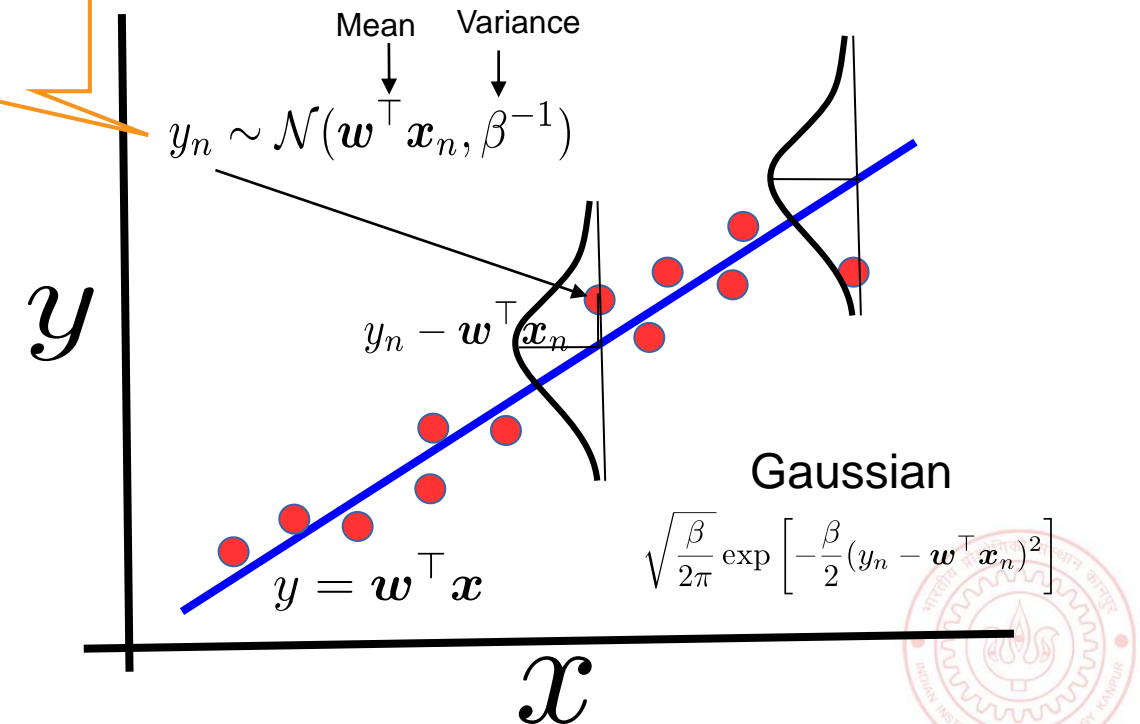
$$y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$$

where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$

Output y_n assumed generated from a Gaussian with mean $\mathbf{w}^\top \mathbf{x}_n$

Each weight assumed real-valued

- Precision (β) variance of the Gaussian noise tells is how noisy the outputs are (i.e., how far from the mean they are)
- Other noise models also possible (e.g., Laplace distribution for noise)



Probabilistic Linear Regression

3

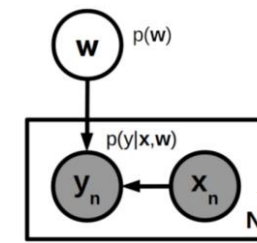


Plate diagram. Hyperparams (λ, β) are fixed and not shown for brevity

- The linear model with Gaussian noise corresponds to a Gaussian likelihood

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

NLL corresponds to squared loss prop. to $(y_n - \mathbf{w}^\top \mathbf{x}_n)^2$

- Assuming responses to be i.i.d. given features and weights

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^\top \mathbf{x}_n, \beta^{-1}) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \beta^{-1} \mathbf{I}_N)$$

$N \times D$ feature matrix

$N \times 1$ response vector

- The above is equivalent to the following

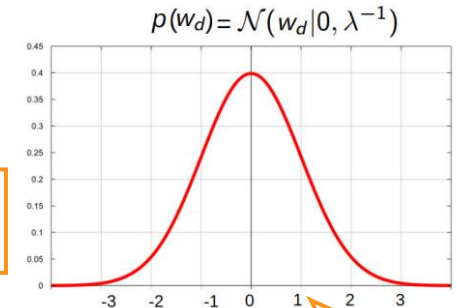
$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon} \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I}_N)$$

- Assume the following **Gaussian prior** on \mathbf{w} ,

Neg. log-prior corresponds to ℓ_2 regularizer with λ being the reg. constant

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

Can even use different λ 's for different w_d 's. Useful in sparse modeling (later)



The precision λ of the Gaussian prior controls how aggressively the prior pushes the elements towards mean (0)

- Then $\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}$ is simply a **linear Gaussian model**

- Can use all the rules of linear Gaussian models to perform inference/predictions 😊

The Posterior

Will only look at fully Bayesian inference. For MLE/MAP, refer to CS771 slides or book



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

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

Must be a Gaussian due to conjugacy

Marginal likelihood for this regression model. Note that it is conditioned on \mathbf{X} too which is assumed given and not being modeled

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

- Using the “completing the squares” trick (or linear Gaussian model results)

$$p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \beta, \lambda) = \mathcal{N}(\mu_N, \Sigma_N)$$

Note that λ and β can be learned under the probabilistic set-up (though assumed fixed as of now)

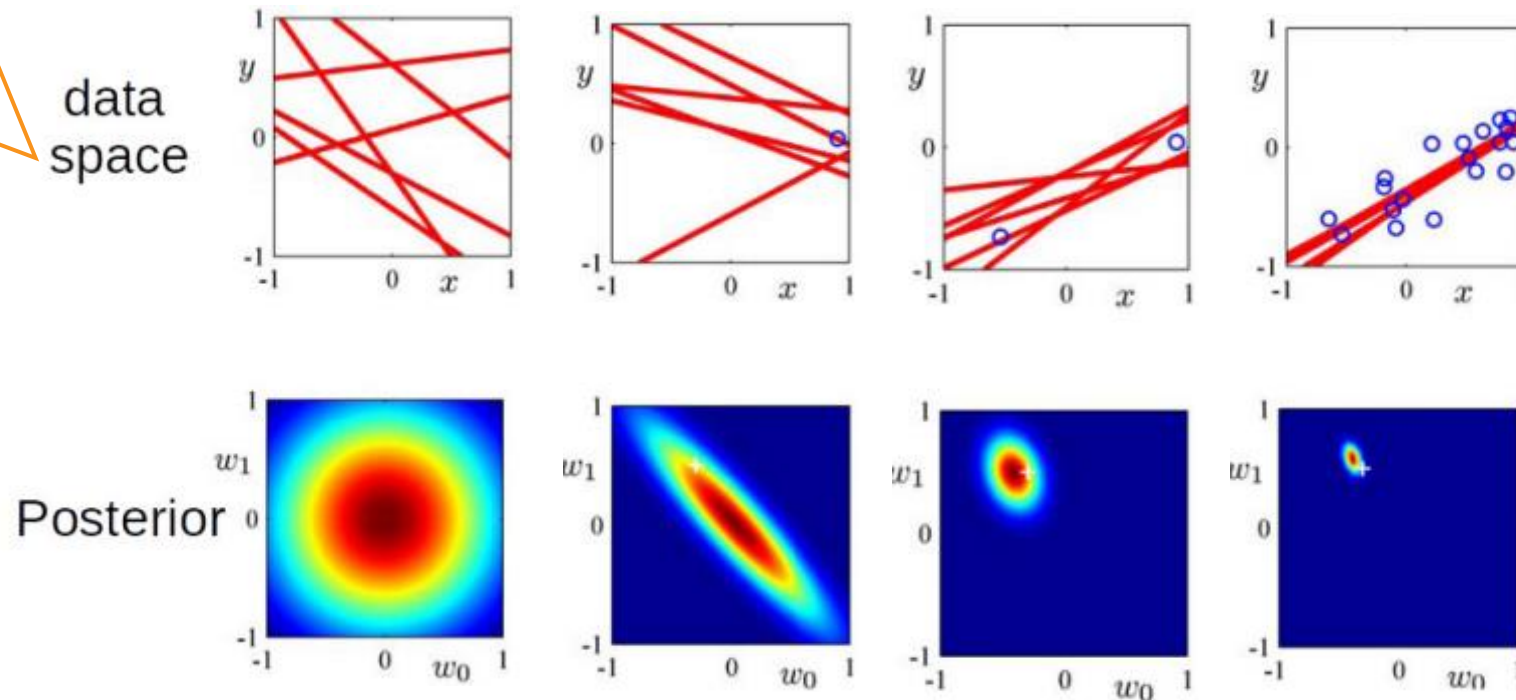
where $\Sigma_N = (\beta \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \lambda \mathbf{I}_D)^{-1} = (\beta \mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1}$ (posterior's covariance matrix)

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

The Posterior: A Visualization

- Assume a lin. reg. problem with true $\mathbf{w} = [w_0, w_1]$, $w_0 = -0.3, w_1 = 0.5$
- 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)

Each red line represents the “data” generated for a randomly drawn \mathbf{w} from the current posterior



Posterior Predictive Distribution

- To get the prediction y_* for a new input \mathbf{x}_* , we can compute its PPD

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

Only \mathbf{w} is unknown with a posterior distribution so only \mathbf{w} has to be integrated out

$\mathcal{N}(y_* | \mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$

$\mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$

- The above is the marginalization of \mathbf{w} from $\mathcal{N}(y_* | \mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$. Using Gaussian results

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

Can also derive it by writing $y_* = \mathbf{w}^\top \mathbf{x}_* + \epsilon$ where $\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$ and $\epsilon \sim \mathcal{N}(0, \beta^{-1})$

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

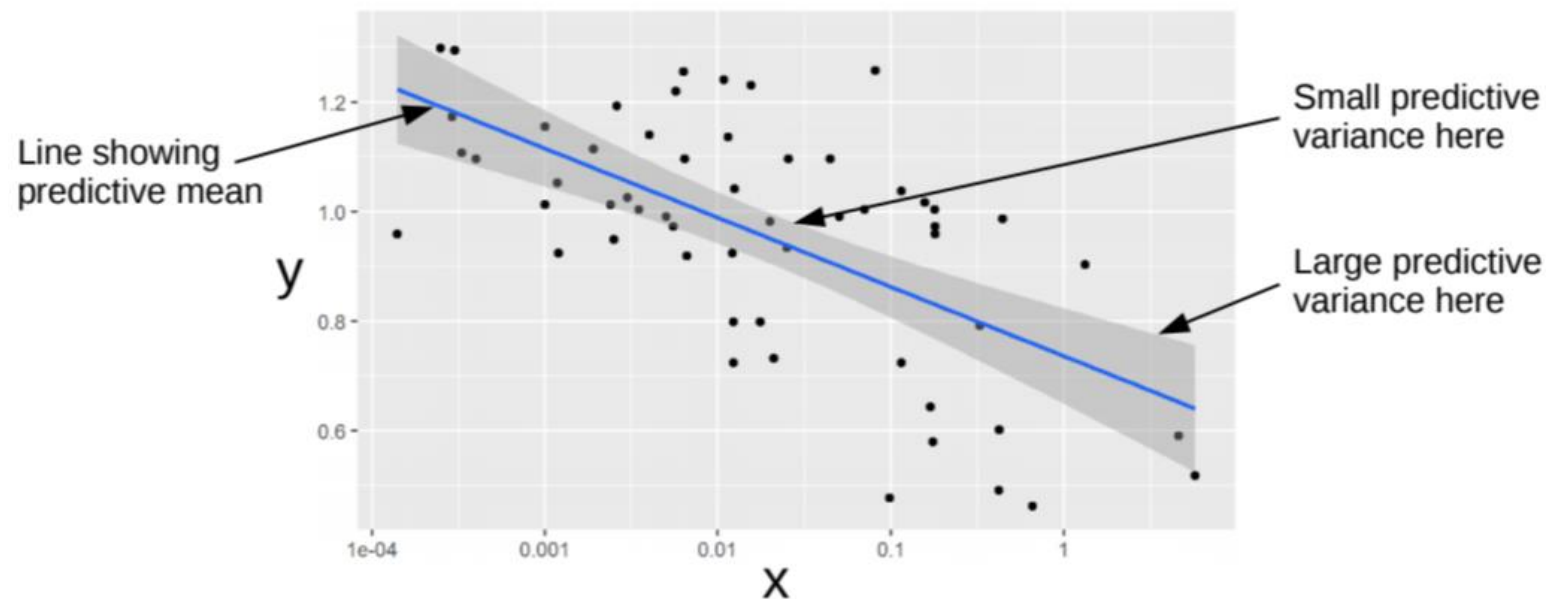
$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \mathcal{N}(\mathbf{w}_{MLE}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MLE prediction} \\ p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \mathcal{N}(\mathbf{w}_{MAP}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MAP prediction} \end{aligned}$$

Since PPD also takes into account the uncertainty in \mathbf{w} , the predictive variance is larger

- Unlike MLE/MAP, 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

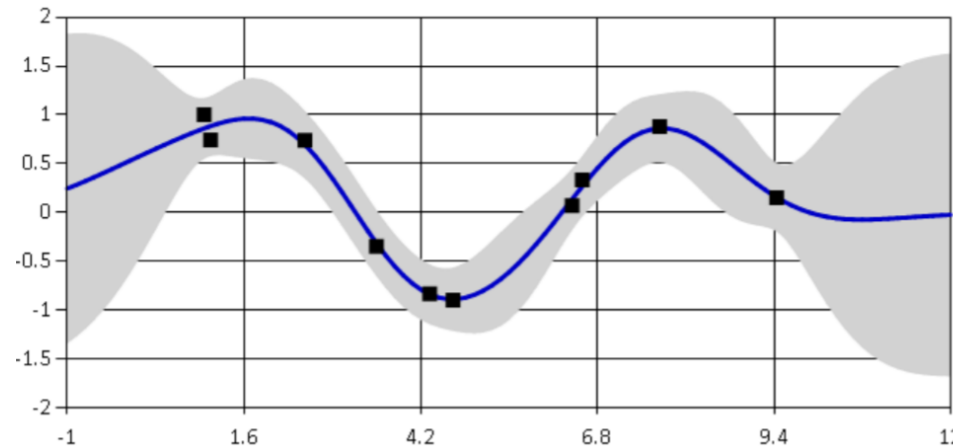
- Black dots are training examples



- 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



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 $\phi(\mathbf{x})$

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

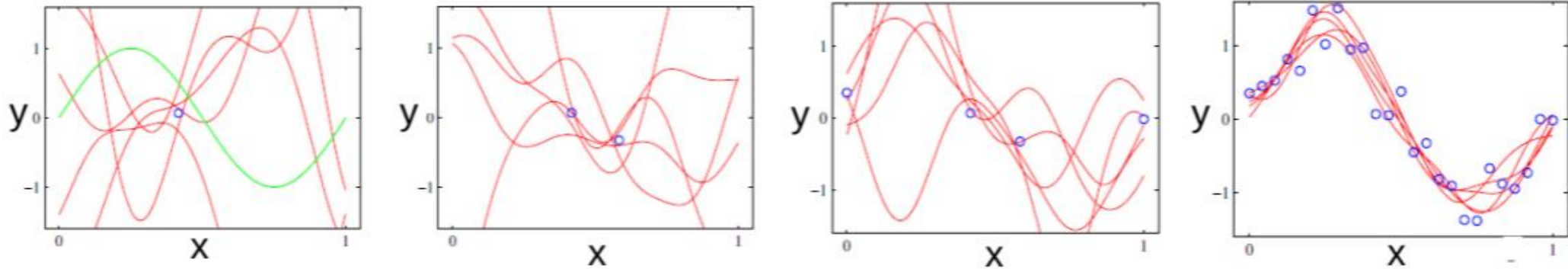
Can be pre-defined or extracted by a pretrained deep neural net

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

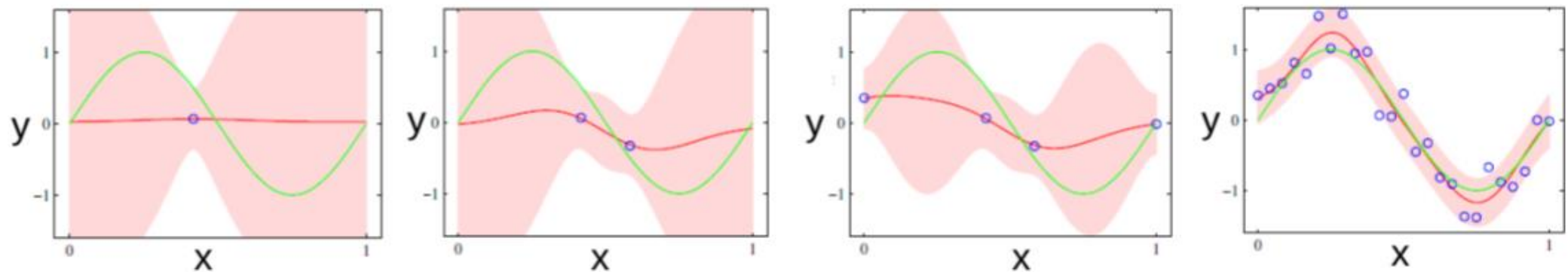


More on Visualization of Uncertainty

- Figures below: Green curve is the true function and blue circles are observations



- Posterior of the nonlinear regression model: Some curves drawn from the posterior

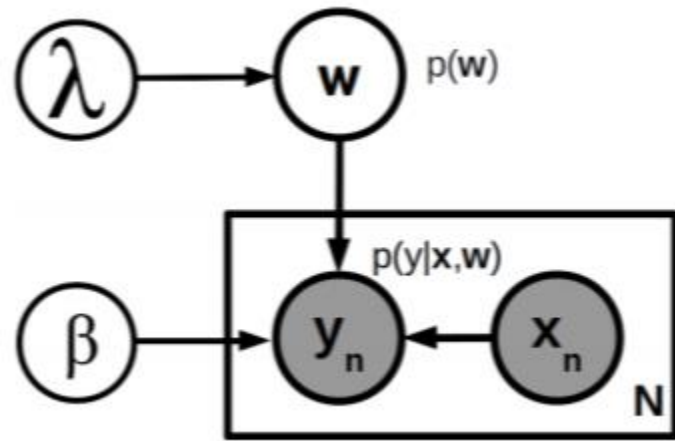


- PPD: Red curve is predictive mean, shaded region denotes predictive uncertainty



Hyperparameters

- The probabilistic linear reg. model we saw had two hyperparams (β, λ)
 - Thus total three unknowns $(\mathbf{w}, \beta, \lambda)$



Need posterior over all the 3 unknowns

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

PPD would require integrating out all 3 unknowns

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

- Posterior and PPD computation is intractable. Several ways to address this
 - MLE-II for (β, λ) : $\hat{\beta}, \hat{\lambda} = \arg \max_{\beta, \lambda} p(\mathbf{y} | \mathbf{X}, \beta, \lambda)$. Use them to infer the posterior of θ and PPD
 - Use alternating estimation like EM (e.g., E step computes θ , M step computes (β, λ))
 - Use MCMC or Variational Inference to approximate the above posterior and PPD



For any model where hyperparams are estimated by MLE-II, the posterior and PPD is approximated in a similar fashion



- For the probabilistic linear regression model, the overall posterior over unknowns

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

- With MLE-II approx of (β, λ) , $p(\beta, \lambda | \mathbf{X}, \mathbf{y}) \approx \delta(\hat{\beta}, \hat{\lambda})$, a point mass at $\hat{\beta}, \hat{\lambda}$

$$p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) \approx p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda})$$

Same as the posterior of \mathbf{w} with the hyperparameters fixed

- Likewise, the PPD will be approximated as follows

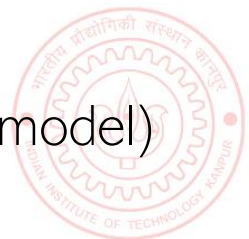
$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}, \beta, \lambda | \mathbf{X}, \mathbf{y}) d\mathbf{w} d\beta d\lambda \\ &= \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \beta, \lambda) p(\beta, \lambda | \mathbf{X}, \mathbf{y}) d\beta d\lambda d\mathbf{w} \\ &\approx \int p(y_* | \mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \hat{\beta}, \hat{\lambda}) d\mathbf{w} \end{aligned}$$

Same form for the PPD as in the case of fixed hyperparams

Only need to integrate over \mathbf{w} , since other two are fixed at their MLE-II solutions

Road So Far and What Lies Ahead

- Seen Bayesian inference for several models with a single unknown parameter (and another simple case where we had two unknown parameters - Gaussian with unknown mean and precision)
- Focused on the cases where the likelihood and prior are conjugate
- Both posterior as well as posterior predictive are computable easily in such cases
- Saw various nice properties of exp. family distributions and parameter estimation for such distributions. Also saw estimation in a conditional model (linear regression)
- Things become more challenging/interesting for more complex models, e.g.,
 - Multiple unknown parameters (e.g., hyperparams, latent variables, hierarchical models etc)
 - Likelihood and prior are not conjugate. Approximate inference methods (MCMC, VI, etc)
- Basic ideas we saw will turn out to be useful in more complex models as well
 - Conditionally-conjugate or locally-conjugate models (conjugacy exists in sub-parts of the model)
 - Some approximate inference methods, e.g., Gibbs sampling, VI, also rely on conjugacy



Coming Up Next

- Hyperparameter estimation for Bayesian linear regression
- Sparse modeling

