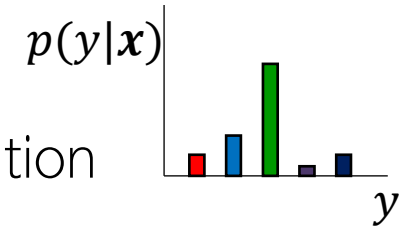# Probabilistic Supervised Learning: Linear Regression

CS772A: Probabilistic Machine Learning

Piyush Rai

# Probabilistic Supervised Learning

- Goal: To learn the conditional distribution $p(y|x)$ of output given input

- The form of the distribution $p(y|x)$ depends on output type, e.g.,

  - Real: Model $p(y|\boldsymbol{x})$ using a Gaussian (or some other suitable real-valued distribution)

  - Binary: Model $p(y|\boldsymbol{x})$ using a Bernoulli

  - Categorical/multiclass: Model $p(y|\boldsymbol{x})$ using a multinoulli/categorical distribution

    $p(y|\boldsymbol{x})$

    $y$

  - Various other types (e.g., count, positive reals, etc) can also be modeled using appropriate distributions (e.g., Poisson for count, gamma for positive reals)

- The distribution $p(y|\boldsymbol{x})$ can be defined directly or indirectly

"Indirect" way requires first learning the joint distribution of inputs and outputs

"Direct" way without modeling the inputs $\boldsymbol{x}_n$

Parameters of this distribution are the outputs of function $f$

"Indirect" way by modeling the outputs as well as the inputs

$$p(y|\boldsymbol{x}) = p(y|f(\boldsymbol{x}, \boldsymbol{w}))$$

$$p(y|\boldsymbol{x}) = \frac{p(y, \boldsymbol{x})}{p(\boldsymbol{x})}$$

# Discriminative vs Generative Sup. Learning

- Direct way of sup. learning is discriminative, indirect way is generative

## Discriminative Approach

$$p(y|\boldsymbol{x}) = p(y|f(\boldsymbol{x}, \boldsymbol{w}))$$

$f$ can be any function which uses inputs and weights $\boldsymbol{w}$ to defines parameters of distr. $p$

### Some examples

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

$$p(y|\boldsymbol{x}) = \text{Bernoulli}(y|\sigma(\boldsymbol{w}^\top \boldsymbol{x}))$$

## Generative Approach

$$p(y|\boldsymbol{x}) = \frac{p(y, \boldsymbol{x})}{p(\boldsymbol{x})}$$

Requires estimating the joint distribution of inputs and outputs to get the conditional $p(y|\boldsymbol{x})$ (unlike the discriminative approach which directly estimates the conditional $p(y|\boldsymbol{x})$ and does not model the distribution of $\boldsymbol{x}$)

- Note: Generative approach can also be used for other settings too, such as unsupervised learning and semi-supervised learning (will see later)

# Probabilistic Linear Regression

A discriminative model for regression problems

- Assume training data $\{x_n, y_n\}_{n=1}^N$, with features $x_n \in \mathbb{R}^D$ and responses $y_n \in \mathbb{R}$

Unknown to be estimated

Each weight assumed real-valued

- Assume $y_n$ generated by a noisy linear model with wts $w = [w_1, \ldots, w_D] \in \mathbb{R}^D$
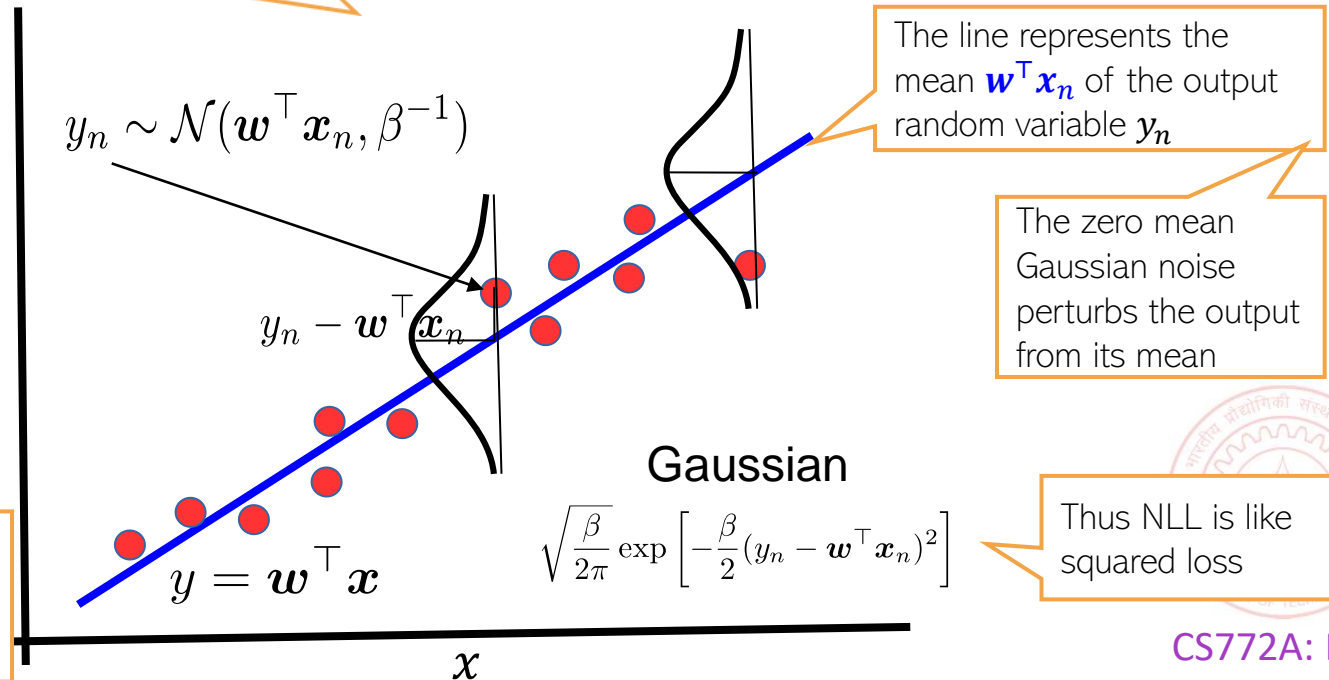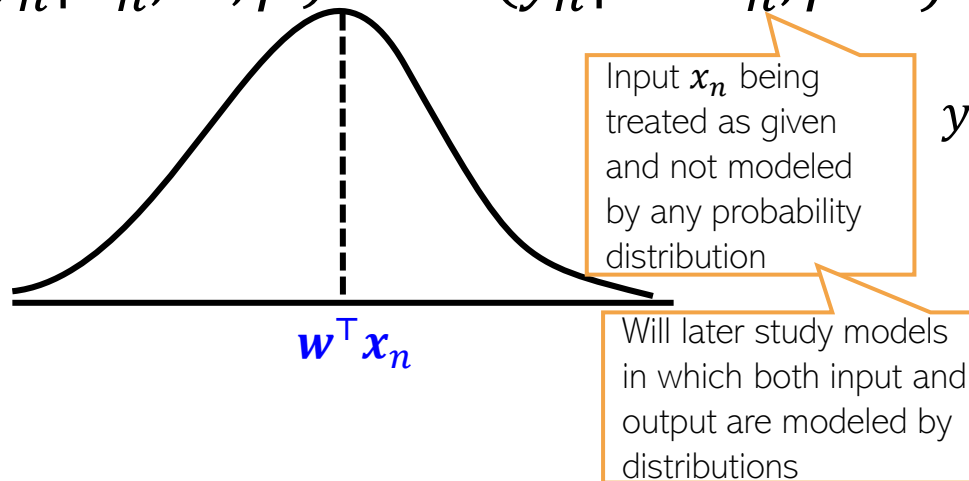
$$y_n = w^\top x_n + \epsilon_n$$

Gaussian noise drawn from $\mathcal{N}(\epsilon_n | 0, \beta^{-1})$

- Notation alert: $\beta$ is the precision of Gaussian noise (and $\beta^{-1}$ the variance)

Unknown to be estimated

Likelihood model

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

$w^\top x_n$

Input $x_n$ being treated as given and not modeled by any probability distribution

Will later study models in which both input and output are modeled by distributions

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

$y_n - w^\top x_n$

$y$

$y = w^\top x$

The line represents the mean $w^\top x_n$ of the output random variable $y_n$

The zero mean Gaussian noise perturbs the output from its mean

Gaussian

$$\sqrt{\frac{\beta}{2\pi}} \exp\left[-\frac{\beta}{2}(y_n - w^\top x_n)^2\right]$$

Thus NLL is like squared loss

$x$

# Probabilistic Linear Regression

- For all the training data, we can write the above model in matrix-vector notation

$\boldsymbol{y} = [y_1; y_2; \ldots; y_N]$ is the $N \times 1$ response vector

$\boldsymbol{X} = [\boldsymbol{x}_1^\mathsf{T}; \boldsymbol{x}_2^\mathsf{T}; \ldots; \boldsymbol{x}_N^\mathsf{T}]$ is the $N \times D$ input matrix
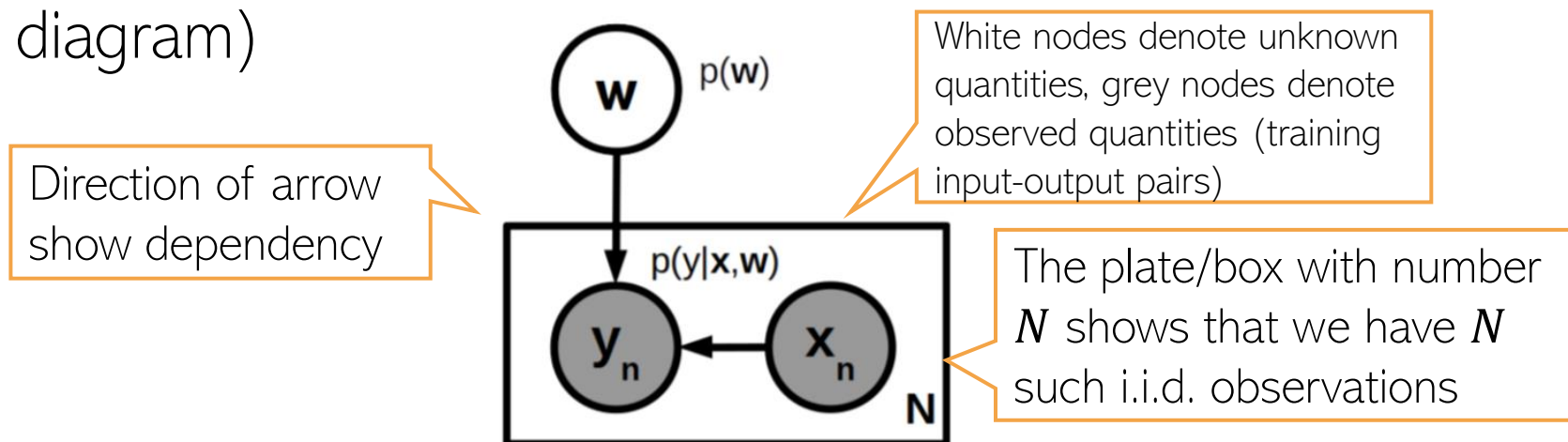
$\boldsymbol{\epsilon} = [\epsilon_1; \epsilon_2; \ldots; \epsilon_N]$ is the $N \times 1$ noise vector drawn from $\mathcal{N}(\boldsymbol{0}, \beta^{-1}\boldsymbol{I}_N)$

Same as writing
$p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w}, \beta) = \mathcal{N}(\boldsymbol{y}|\boldsymbol{X}\boldsymbol{w}, \beta^{-1}\boldsymbol{I}_N)$

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{w} + \boldsymbol{\epsilon}$$

- This is a linear Gaussian model with $\boldsymbol{w}$ being the unknown Gaussian r.v.

- A simple "plate diagram" for this model would look like this (hyperparameters not shown in the diagram)

White nodes denote unknown quantities, grey nodes denote observed quantities (training input-output pairs)

Direction of arrow show dependency

The plate/box with number $N$ shows that we have $N$ such i.i.d. observations

w p(w)

p(y|x,w)

y_n ← x_n

N

# On compact notations..

- When writing the likelihood (assuming $y_n$'s are i.i.d. given $\boldsymbol{w}$ and $\boldsymbol{x_n}$)

$$p(\boldsymbol{y}|\boldsymbol{X}, \boldsymbol{w}, \beta) = \prod_{n=1}^{N} \mathcal{N}(y_n|\boldsymbol{w}^\top \boldsymbol{x_n}, \beta^{-1})$$

$$= \mathcal{N}(\boldsymbol{y}|\boldsymbol{X}\boldsymbol{w}, \beta^{-1}\boldsymbol{I}_N)$$

- Thus a product of $N$ univariate Gaussians here (not always) is equivalent to an $N$-dim Gaussian over the vector $\boldsymbol{y} = [y_1, y_2, \dots, y_N]$

- We will prefer to use this equivalence at other places too whenever we have multiple i.i.d. random variables, each having a univariate Gaussian distribution
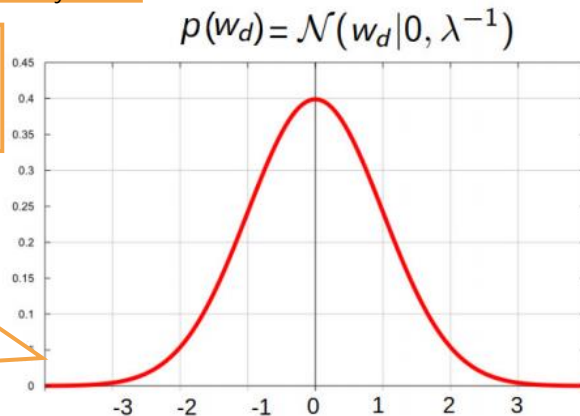
# Prior on weights

- Assume a <span style="color:red">zero-mean Gaussian prior</span> on $\boldsymbol{w}$

> May also use a non-zero mean Gaussian prior, e.g., $\mathcal{N}(w_d|\mu, \lambda^{-1})$ if we expect weights to be close to some value $\mu$

> This prior assumes that *a priori* each weight has a small value (close to zero)

$$p(\boldsymbol{w}|\lambda) = \prod_{d=1}^{D} p(w_d|\lambda) = \prod_{d} \mathcal{N}(w_d|0, \lambda^{-1})$$

> In zero-mean case, $\lambda$ sort of denotes each feature's importance. Think why?

> $\lambda$ controls the uncertainty around our prior belief about value of $w_d$

> Can also use a <span style="color:blue">full covariance matrix $\Lambda^{-1}$</span> for the prior to impose a priori correlations among different weights

> Large $\lambda$ means more aggressive push towards zero

$p(w_d) = \mathcal{N}(w_d|0, \lambda^{-1})$

$$= \mathcal{N}(\boldsymbol{w}|\boldsymbol{0}, \lambda^{-1}\mathbf{I}_D)$$

> The precision $\lambda$ controls how aggressively the prior pushes $w_d$ towards mean (0)

> <span style="color:green">Prior's hyperparameters $(\lambda/\Lambda/\mu)$ etc can be learned as well using point estimation (e.g., MLE-II) or fully Bayesian inference</span>

$$\propto \left(\frac{\lambda}{2\pi}\right)^{\frac{D}{2}} \exp\left[-\frac{\lambda}{2}\, \boldsymbol{w}^{\top}\boldsymbol{w}\right]$$

- Zero-mean Gaussian prior corresponds to $\ell_2$ regularizer

> Reason: The negative log prior $-\log p(\boldsymbol{w}) \propto \frac{\lambda}{2}\, \boldsymbol{w}^{\top}\boldsymbol{w}$

# The Posterior

MLE/MAP left as an exercise

- The posterior over $w$ (for now, assume hyperparams $\beta$ and $\lambda$ to be known)

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

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

Must be a Gaussian due to conjugacy

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

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

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

Note that $\lambda$ and $\beta$ can be learned under the probabilistic set-up (though assumed fixed as of now)

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

The form is also similar to the solution to ridge regression

$$\text{argmin}_w ||y - Xw||^2 + \lambda w^\top w = (X^\top X + \lambda I)^{-1} X^\top y$$

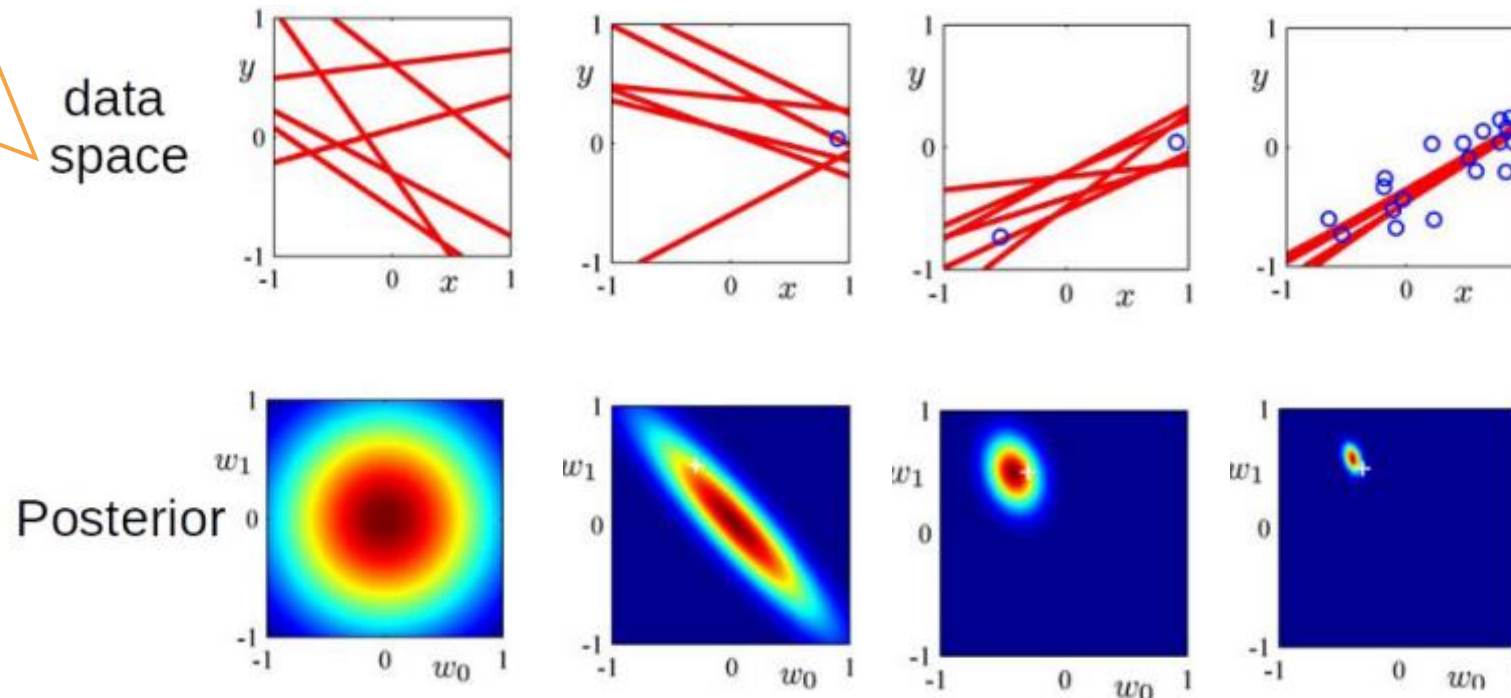MAP solution turns out to be exactly the same (reason: Gaussian's mean and mode are the same)

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

: PML

# The Posterior: A Visualization

- Assume a lin. reg. problem with true $\boldsymbol{w} = [w_0, w_1], w_0 = -0.3, w_1 = 0.5$
- Assume data generated by a linear regression model $y = w_0 + w_1 x +$ "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 $\boldsymbol{w}$ for different number of observations (note: with no observations, the posterior = prior)

Each red line represents the "data" generated for a randomly drawn $\boldsymbol{w}$ from the current posterior

data space

Posterior

# Posterior Predictive Distribution

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

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

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

$$\mathcal{N}(y_*|\boldsymbol{w}^\top \boldsymbol{x}_*, \beta^{-1}) \qquad \mathcal{N}(\boldsymbol{w}|\boldsymbol{\mu}_N, \boldsymbol{\Sigma}_N)$$

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

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

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

- So we have a predictive mean $\boldsymbol{\mu}_N^\top \boldsymbol{x}_*$ as well as an input-specific predictive variance

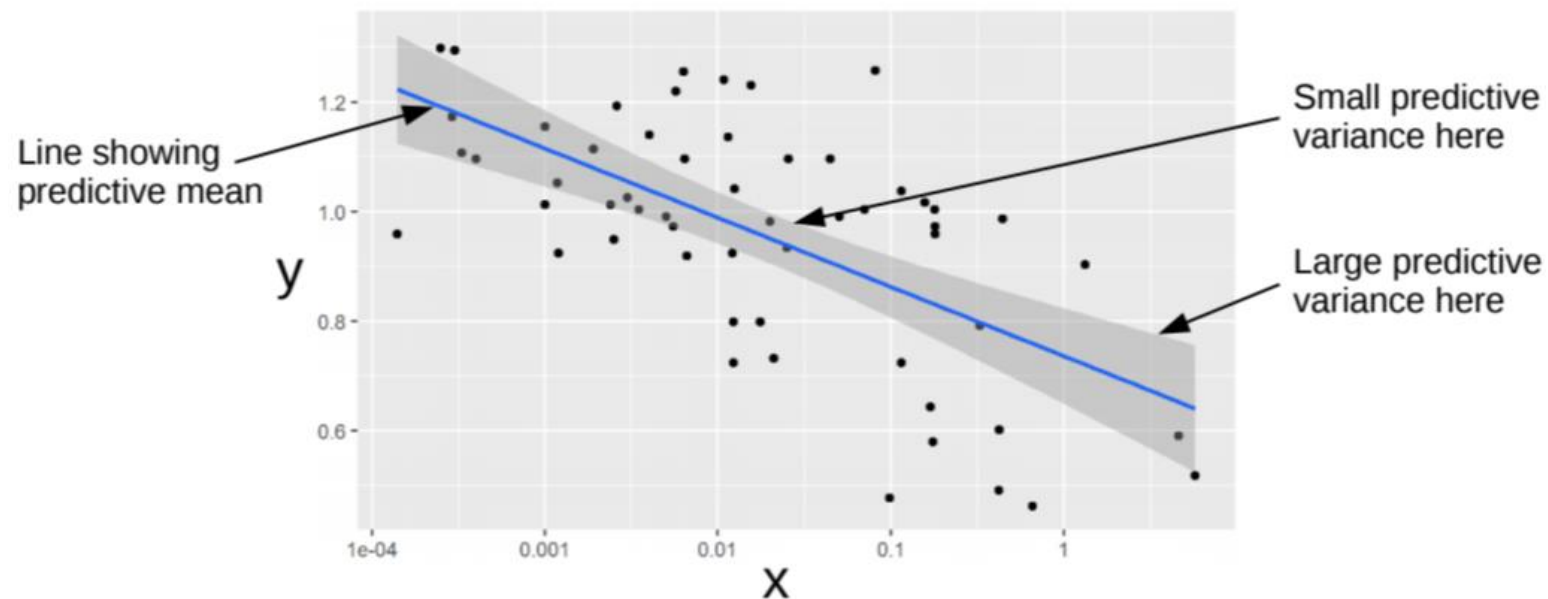- In contrast, MLE and MAP make "plug-in" predictions (using the point estimate of $\boldsymbol{w}$)

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

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

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

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

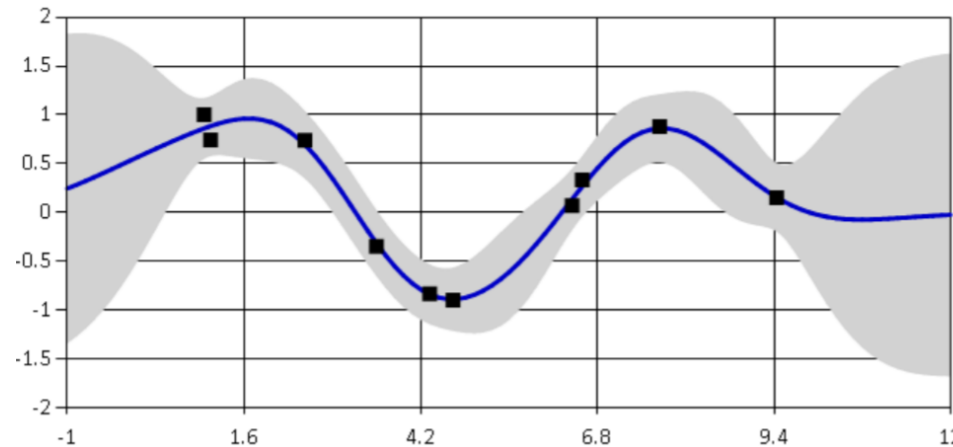# Posterior Predictive Distribution: An Illustration

- Black dots are training examples



- Width of the shaded region at any $x$ denotes the predictive uncertainty at that $x$ (+/- 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 $\boldsymbol{x}$ by a nonlinear mapping $\boldsymbol{\phi(x)}$

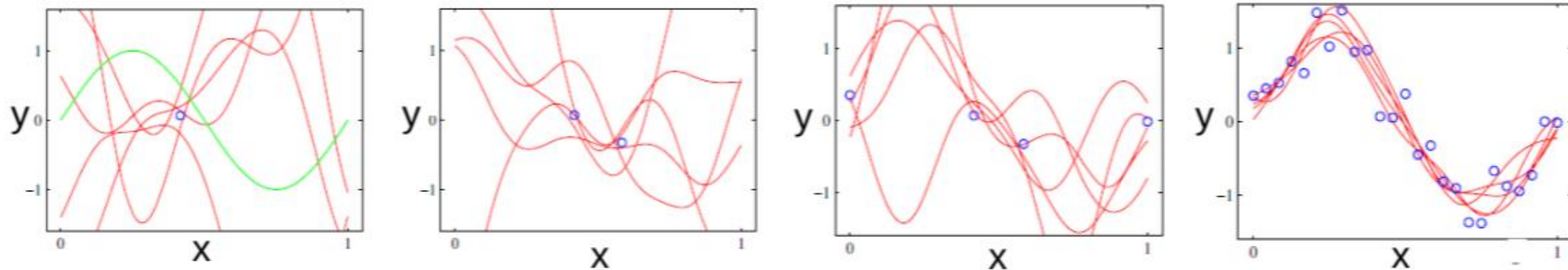$$p(y|\boldsymbol{x}, \boldsymbol{w}) = \mathcal{N}(\boldsymbol{w}^\top \phi(\boldsymbol{x}), \beta^{-1})$$

> Can be pre-defined (e.g., replace a scalar $x$ by polynomial mapping $[1, x, x^2]$) or extracted by a pretrained deep neural net

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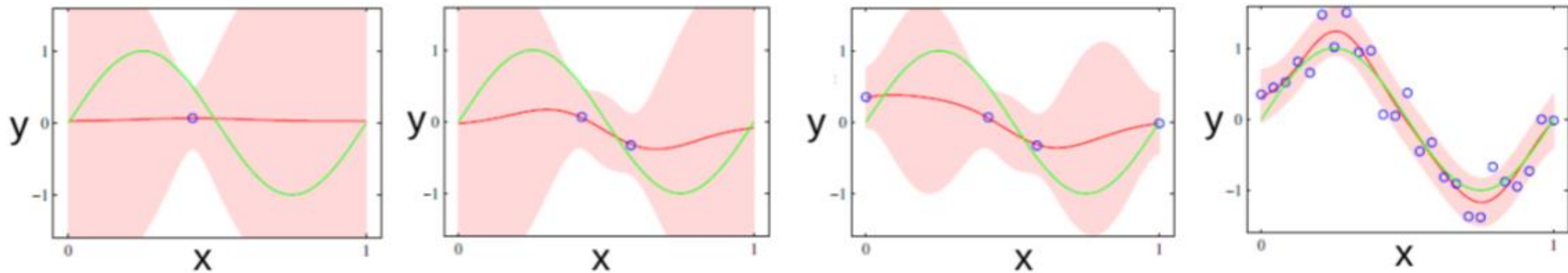- More on nonlinear regression when we discuss Gaussian Processes

# More on Visualization of Uncertainty

- Figures below: Green curve is the true function and blue circles are observations
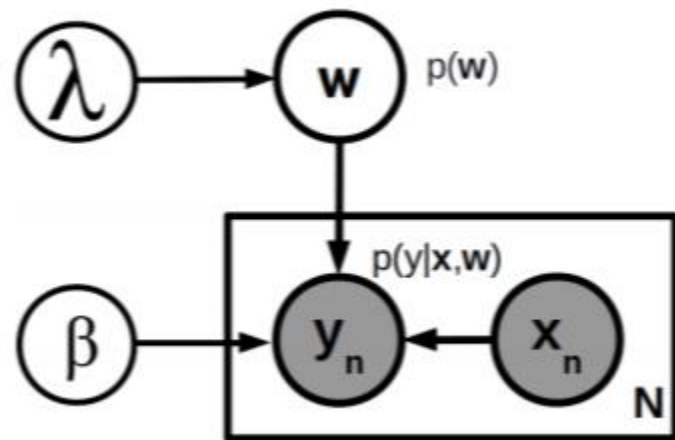- Posterior of the nonlinear regression model: Some curves drawn from the posterior



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

# Estimating Hyperparameters via MLE-II

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

Need posterior over all the 3 unknowns

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

PPD would require integrating out all 3 unknowns

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

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

- Posterior and PPD computation is intractable.

Called "MLE-II" because we are maximizing marginal likelihood, not the likelihood

- If we just want point estimates for $(\beta, \lambda)$ then MLE-II is an option

Will see various other methods like EM, variational inference, MCMC, etc later

And then compute $p(\boldsymbol{w}|\boldsymbol{X}, \boldsymbol{y}, \hat{\beta}, \hat{\lambda})$ treating $\hat{\beta}, \hat{\lambda}$ as given

$$(\hat{\beta}, \hat{\lambda}) = \operatorname{argmax}_{\beta, \lambda} \log p(\boldsymbol{y}|\boldsymbol{X}, \beta, \lambda)$$

For regression with Gaussian likelihood and Gaussian prior on $\boldsymbol{w}$, the marginal likelihood has an exact expression

# Prob. Linear Regression: Some Other Variations

- Can use other likelihoods $p(y_n|\boldsymbol{x}_n, \boldsymbol{w})$ and/or prior distribution $p(\boldsymbol{w})$

- Laplace distribution for the likelihood

$$p(y_n|\boldsymbol{x}_n, \boldsymbol{w}) = \text{Lap}(y_n|\boldsymbol{w}^\top \boldsymbol{x}_n, b)$$

- Heteroskedastic noise in the likelihood, e.g.,

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

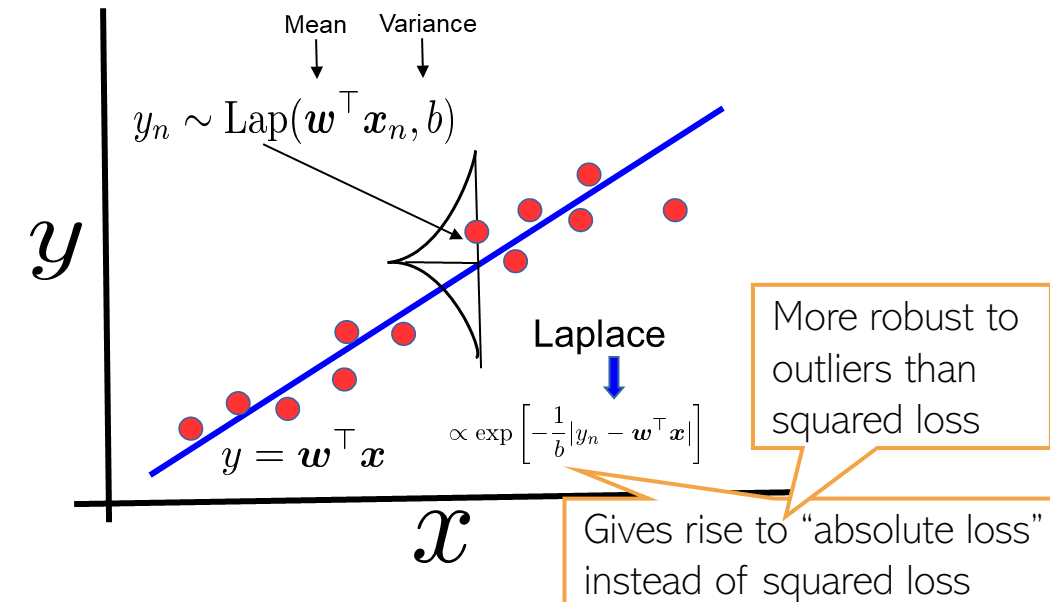> Can even assume $\beta_n$ to depend on input $\boldsymbol{x}_n$

> Different noise distribution $\mathcal{N}(0, \beta_n^{-1})$ for each $y_n$

- Feature-specific variances in the prior for $\boldsymbol{w}$

$$p(\boldsymbol{w}) = \prod_{d=1}^{D} \mathcal{N}(w_d|0, \lambda_d^{-1}) = \mathcal{N}(\boldsymbol{w}|\boldsymbol{0}, \boldsymbol{\Lambda}^{-1})$$

> This has the effect of having feature-specific regularization

> Diagonal precision/covariance matrix with $\lambda_d$'s along the columns of $\Lambda$

> Since we can also learn these precisions (e.g., using MLE-II), using such a prior, we can learn the importance of different features (feature selection) which isn't possible with a $\mathcal{N}(\boldsymbol{w}|\boldsymbol{0}, \lambda^{-1}\boldsymbol{I})$ prior with spherical covariance

Mean    Variance

$y_n \sim \text{Lap}(\boldsymbol{w}^\top \boldsymbol{x}_n, b)$

$y$

Laplace

$y = \boldsymbol{w}^\top \boldsymbol{x}$    $\propto \exp\left[-\frac{1}{b}|y_n - \boldsymbol{w}^\top \boldsymbol{x}|\right]$

$x$

> More robust to outliers than squared loss

> Gives rise to "absolute loss" instead of squared loss