

Gaussian Observation Model: Some Examples

CS772A: Probabilistic Machine Learning

Piyush Rai

Plan today

- Two models with observations as Gaussian distributed (i.e., Gaussian likelihood)

Estimate the mean μ of this Gaussian given N i.i.d. training observations

$$p(y_n|\mu, \sigma^2) = \mathcal{N}(y_n|\mu, \sigma^2)$$

Estimate the weight vector \mathbf{w} of this (probabilistic) linear regression model given N i.i.d. training observations

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

- We will look at computation of the posterior distribution (MLE/MAP left as an exercise) over the parameters as well as the predictive distributions



Recap: Linear Gaussian Model (LGM)

- LGM defines a noisy **lin. transform** of a Gaussian r.v. $\boldsymbol{\theta}$ with $p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Both $\boldsymbol{\theta}$ and \mathbf{y} are vectors (can be of different sizes)

Also assume $\mathbf{A}, \mathbf{b}, \boldsymbol{\Lambda}, \mathbf{L}$ to be known; only $\boldsymbol{\theta}$ is unknown

$$\mathbf{y} = \mathbf{A}\boldsymbol{\theta} + \mathbf{b} + \boldsymbol{\epsilon}$$

Noise vector - independently and drawn from $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \mathbf{L}^{-1})$

- Easy to see that, conditioned on $\boldsymbol{\theta}$, \mathbf{y} too has a Gaussian distribution

Conditional distribution

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

- Assume $p(\boldsymbol{\theta})$ as prior and $p(\mathbf{y}|\boldsymbol{\theta})$ as the likelihood, and defining $\boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^\top \mathbf{L} \mathbf{A})^{-1}$

Posterior of $\boldsymbol{\theta}$

$$p(\boldsymbol{\theta}|\mathbf{y}) = \frac{p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\mathbf{y})} = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\Sigma}(\mathbf{A}^\top \mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}), \boldsymbol{\Sigma})$$

Marginal distribution

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

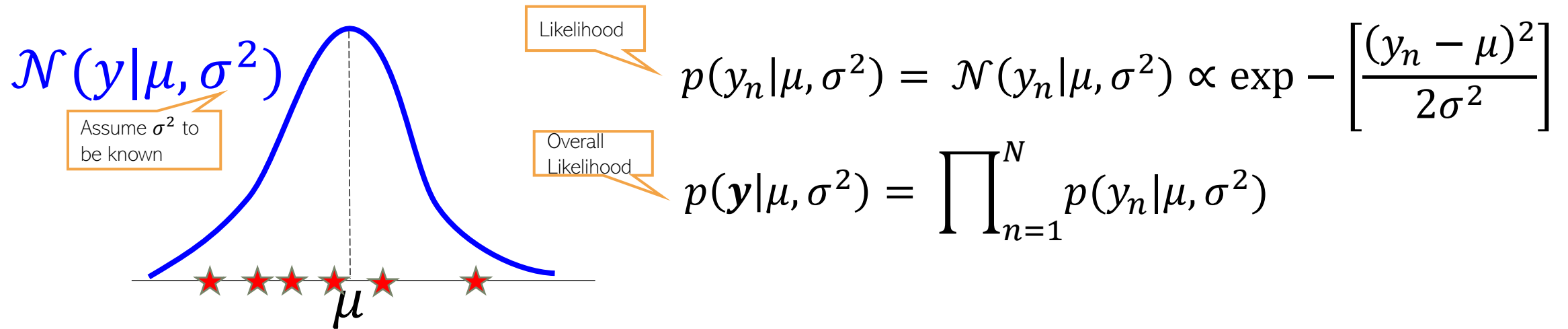
- Many probabilistic ML models are LGMs
- These results are very widely used (PRML Chap. 2 contains a proof)



Posterior Distribution for Gaussian's Mean

Its MLE/MAP estimation left as an exercise

- Given: N i.i.d. scalar observations $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$ assumed drawn from $\mathcal{N}(y|\mu, \sigma^2)$



- Note: Easy to see that each y_n drawn from $\mathcal{N}(y|\mu, \sigma^2)$ is equivalent to the following

Thus y_n is like a noisy version of μ with zero mean Gaussian noise added to it

$$y_n = \mu + \epsilon_n \quad \text{where } \epsilon_n \sim \mathcal{N}(0, \sigma^2)$$

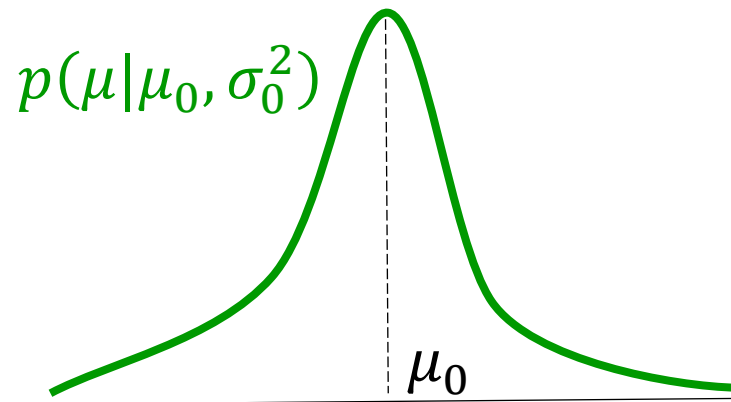
- Let's estimate mean μ given \mathbf{y} using fully Bayesian inference (not point estimation)



A prior distribution for the mean

- To compute posterior, need a prior over μ
- Let's choose a Gaussian prior

$$p(\mu|\mu_0, \sigma_0^2) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) \\ \propto \exp \left[-\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$



- The prior basically says that a priori we believe μ is close to μ_0
- The prior's variance σ_0^2 denotes how certain we are about our belief
- We will assume that the prior's hyperparameters (μ_0, σ_0^2) are known
- Since σ^2 in the likelihood $\mathcal{N}(y|\mu, \sigma^2)$ is known, Gaussian prior $\mathcal{N}(\mu|\mu_0, \sigma_0^2)$ on μ is also conjugate to the likelihood (thus posterior of μ will also be Gaussian)



The posterior distribution for the mean

- The posterior distribution for the unknown mean parameter μ

On conditioning side, skipping all fixed params and hyperparams from the notation

$$p(\mu|\mathbf{y}) = \frac{p(\mathbf{y}|\mu)p(\mu)}{p(\mathbf{y})} \propto \prod_{n=1}^N \exp \left[-\frac{(y_n - \mu)^2}{2\sigma^2} \right] \exp \left[-\frac{(\mu - \mu_0)^2}{2\sigma_0^2} \right]$$

- Easy to see that the above will be prop. to **exp of a quadratic function** of μ . Simplifying:

$$p(\mu|\mathbf{y}) \propto \exp \left[-\frac{(\mu - \mu_N)^2}{2\sigma_N^2} \right]$$

Gaussian posterior's precision is the sum of the prior's precision and sum of the noise precisions of all the observations

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}$$

Contribution from the prior

Contribution from the data

Also the MLE solution for μ

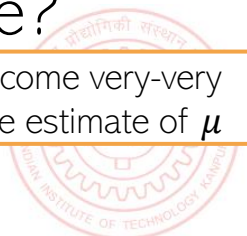
Gaussian posterior's mean is a convex combination of prior's mean and the MLE solution

$$\mu_N = \frac{N\sigma_0^2}{N\sigma_0^2 + \sigma^2} \mu_0 + \frac{\sigma^2}{N\sigma_0^2 + \sigma^2} \bar{y} \quad (\text{where } \bar{y} = \frac{\sum_{n=1}^N y_n}{N})$$

- What happens to the posterior as N (number of observations) grows very large?

- Data (likelihood part) overwhelms the prior
- Posterior's variance σ_N^2 will approximately be σ^2/N (and goes to 0 as $N \rightarrow \infty$)
- The posterior's mean μ_N approaches \bar{y} (which is also the MLE solution)

Meaning, we become very-very certain about the estimate of μ



The Predictive Distribution

- If given a point estimate $\hat{\mu}$, the plug-in predictive distribution for a test y_* would be

This is an approximation of the true PPD $p(y_*|\mathbf{y})$

The best point estimate

$$p(y_*|\hat{\mu}, \sigma^2) = \mathcal{N}(y_*|\hat{\mu}, \sigma^2)$$

- On the other hand, the posterior predictive distribution of y_* would be

$$\begin{aligned} p(y_*|\mathbf{y}) &= \int p(y_*|\mu, \sigma^2)p(\mu|\mathbf{y})d\mu \\ &= \int \mathcal{N}(y_*|\mu, \sigma^2)\mathcal{N}(\mu|\mu_N, \sigma_N^2)d\mu \\ &= \mathcal{N}(y_*|\mu_N, \sigma^2 + \sigma_N^2) \end{aligned}$$

This "extra" variance σ_N^2 in PPD is due to the averaging over the posterior's uncertainty

If conditional is Gaussian then marginal is also Gaussian

A useful fact: When we have conjugacy, the posterior predictive distribution also has a closed form (will see this result more formally when talking about exponential family distributions)



PRML [Bis 06], 2.115, and also mentioned in prob-stats refresher slides

- For an alternative way to get the above result, note that, for test data

$$y_* = \mu + \epsilon \quad \mu \sim \mathcal{N}(\mu_N, \sigma_N^2) \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Using the posterior of μ since we are at test stage now

$$\Rightarrow p(y_*|\mathbf{y}) = \mathcal{N}(y_*|\mu_N, \sigma^2 + \sigma_N^2)$$

Since both μ and ϵ are Gaussian r.v., and are independent, y_* also has a Gaussian posterior predictive, and the respective means and variances of μ and ϵ get added up



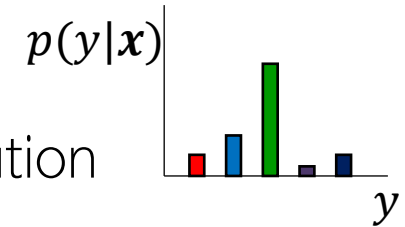
Gaussian Observation Model: Summary/Some Facts

- MLE/MAP for μ, σ^2 (or both) is straightforward in Gaussian observation models.
- Posterior also straightforward in most situations for such models
 - (As we saw) computing posterior of μ is easy (using Gaussian prior) if variance σ^2 is known
 - Likewise, computing posterior of σ^2 is easy (using **gamma prior** on σ^2) if mean μ is known
- If μ, σ^2 both are unknown, posterior computation requires computing $p(\mu, \sigma^2 | \mathbf{y})$
 - Computing joint posterior $p(\mu, \sigma^2 | \mathbf{y})$ exactly requires a jointly conjugate prior $p(\mu, \sigma^2)$
 - “**Gaussian-gamma**” (“Normal-gamma”) is such a conjugate prior – a product of normal and gamma
 - Note: Computing joint posteriors exactly is possible only in rare cases such this one
- If each observation $\mathbf{y}_n \in \mathbb{R}^D$, can assume a likelihood/observation model $\mathcal{N}(\mathbf{y} | \mu, \Sigma)$
 - Need to estimate a **vector-valued** mean $\mu \in \mathbb{R}^D$. Can use a **multivariate Gaussian prior**
 - Need to estimate a $D \times D$ positive definite covariance **matrix** Σ . Can use a **Wishart prior**
 - If μ, Σ both are unknown, can use **Normal-Wishart** as a conjugate prior



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 $p(y|x)$ using a Gaussian (or some other suitable real-valued distribution)
 - Binary: Model $p(y|x)$ using a Bernoulli
 - Categorical/multiclass: Model $p(y|x)$ using a multinoulli/categorical distribution
 - 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 $p(y|x)$ can be defined directly or indirectly



“Direct” way without modeling the inputs x_n

Parameters of this distribution are the outputs of function f

“Indirect” way by modeling the outputs as well as the inputs

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

$$p(y|x) = p(y|f(x, w))$$

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



Discriminative vs Generative Sup. Learning

Non-probabilistic supervised learning approaches (e.g., SVM) are usually considered discriminative since $p(\mathbf{x})$ is never modeled

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

Discriminative Approach

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

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

Some examples

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

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

Generative Approach

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

Requires estimating the **joint distribution** of inputs and outputs to get the conditional $p(y|\mathbf{x})$ (unlike the discriminative approach which directly estimates the conditional $p(y|\mathbf{x})$ and does not model the distribution of \mathbf{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 $\{\mathbf{x}_n, \mathbf{y}_n\}_{n=1}^N$, with features $\mathbf{x}_n \in \mathbb{R}^D$ and responses $\mathbf{y}_n \in \mathbb{R}$
- Assume \mathbf{y}_n generated by a noisy linear model with wts $\mathbf{w} = [w_1, \dots, w_D] \in \mathbb{R}^D$

Unknown to be estimated

Each weight assumed real-valued

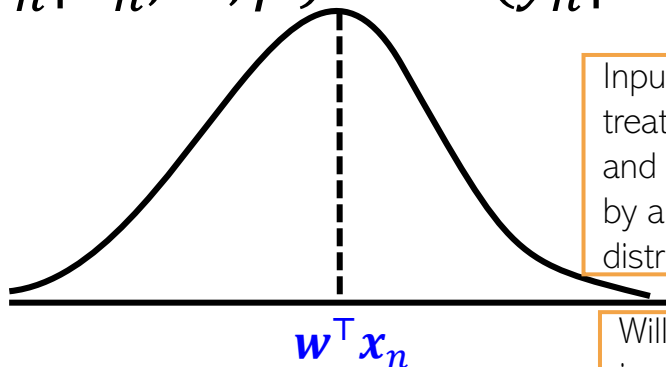
$$\mathbf{y}_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$$

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

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

Likelihood model

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

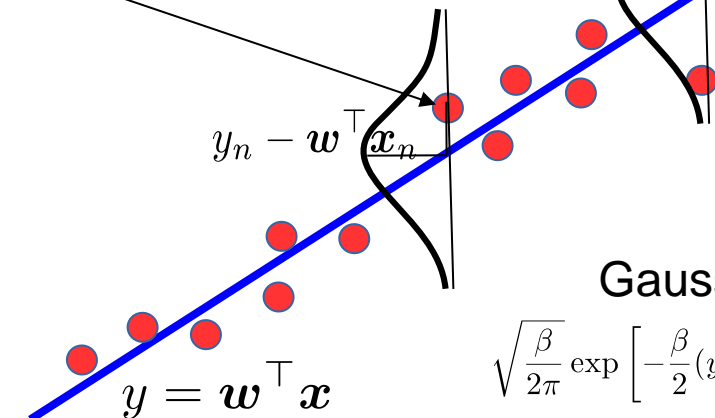


Input \mathbf{x}_n being treated as given and not modeled by any probability distribution

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

Unknown to be estimated

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



The line represents the mean $\mathbf{w}^\top \mathbf{x}_n$ of the output random variable \mathbf{y}_n

The zero mean Gaussian noise perturbs the output from its mean

Gaussian

$$\sqrt{\frac{\beta}{2\pi}} \exp \left[-\frac{\beta}{2} (\mathbf{y}_n - \mathbf{w}^\top \mathbf{x}_n)^2 \right]$$

Thus NLL is like squared loss

Prior on weights

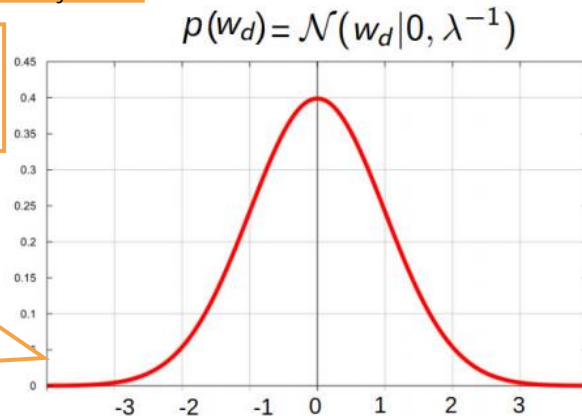
- Assume a **zero-mean Gaussian prior** on \mathbf{w}

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

In zero-mean case, λ sort of denotes each feature's importance. Think why?

Large λ means more aggressive push towards zero

The precision λ controls how aggressively the prior pushes w_d towards mean (0)



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

λ controls the uncertainty around our prior belief about value of w_d

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

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

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 μ



Can also use a **full covariance matrix** Λ^{-1} for the prior to impose a priori correlations among different weights

Prior's hyperparameters ($\lambda/\Lambda/\mu$) etc can be learned as well using point estimation (e.g., MLE-II) or fully Bayesian inference

- Zero-mean Gaussian prior corresponds to ℓ_2 regularizer

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



Probabilistic Linear Regression

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

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

$\mathbf{X} = [\mathbf{x}_1^\top; \mathbf{x}_2^\top; \dots; \mathbf{x}_N^\top]$ is the $N \times D$ input matrix

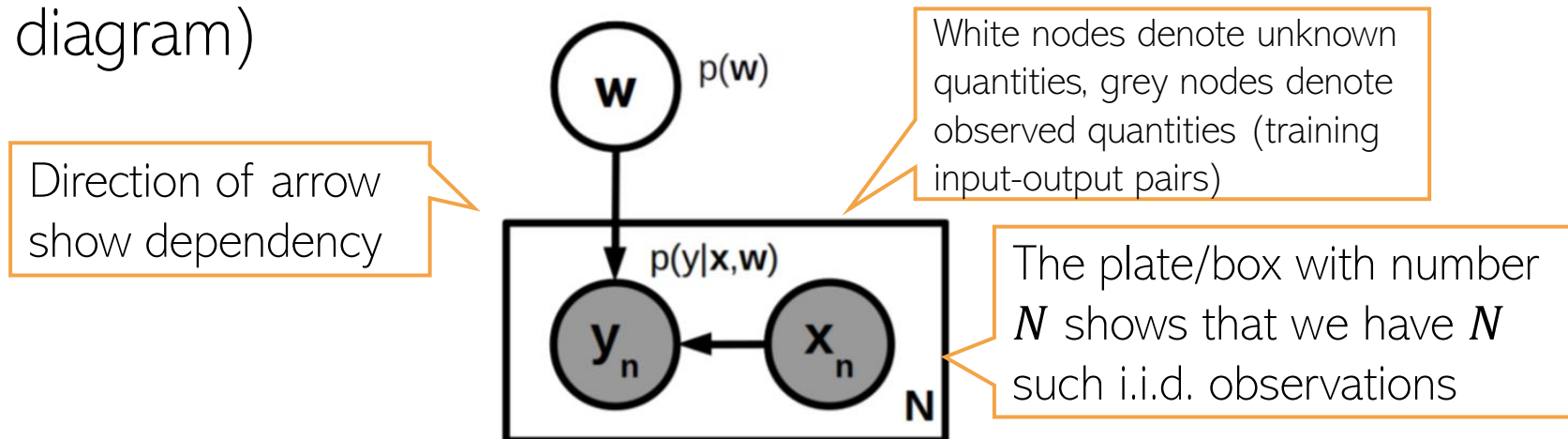
$\boldsymbol{\epsilon} = [\epsilon_1; \epsilon_2; \dots; \epsilon_N]$ is the $N \times 1$ noise vector drawn from $\mathcal{N}(\mathbf{0}, \beta^{-1} \mathbf{I}_N)$

Same as writing

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

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

- Linear Gaussian model and \mathbf{w} is the Gaussian r.v. with $p(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1} \mathbf{I}_D)$
- A simple “plate diagram” for this model would look like this (hyperparameters not shown in the diagram)



On compact notations..

- When writing the likelihood (assuming \mathbf{y}_n 's are i.i.d. given \mathbf{w} and \mathbf{x}_n)

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

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



The Posterior

MLE/MAP left
as an exercise



- 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)

The form is also similar to the solution to ridge regression
 $\text{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 + \lambda \mathbf{w}^\top \mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

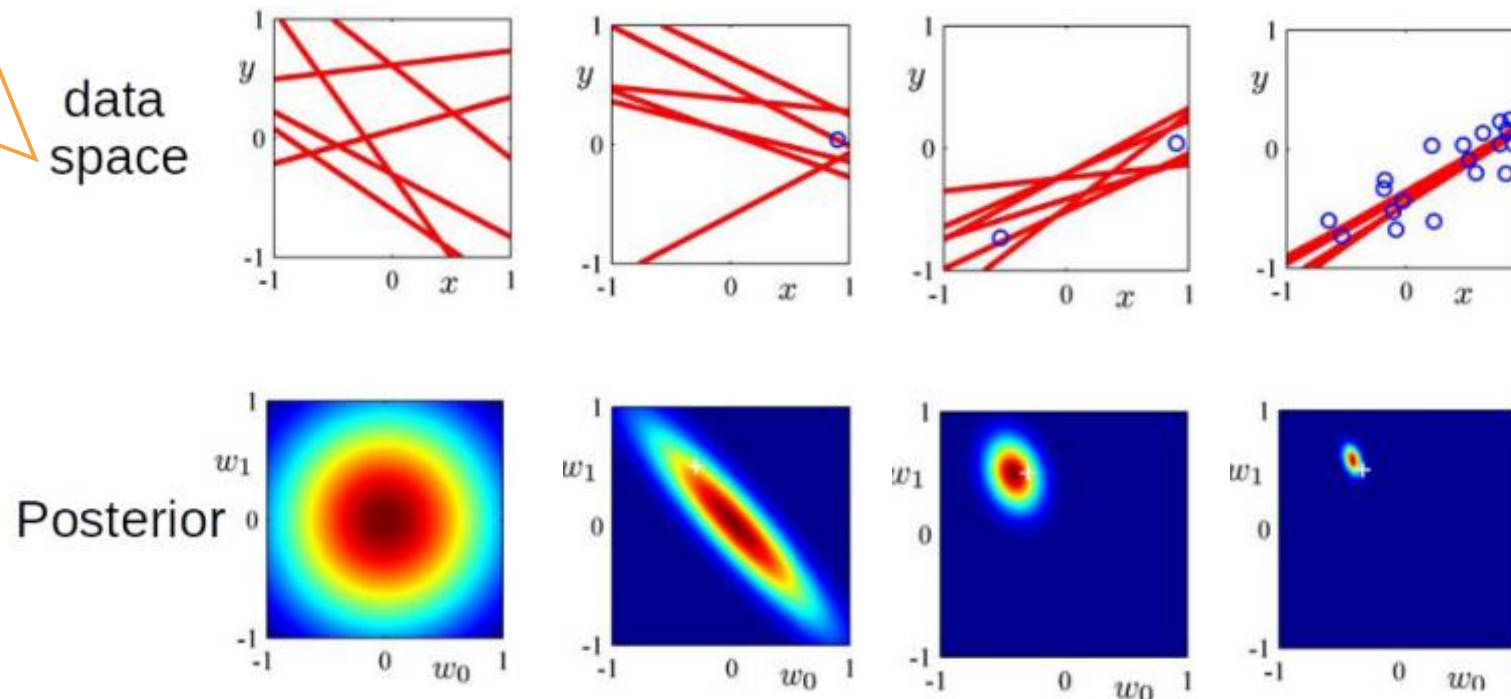
MAP solution turns out to be exactly
the same (reason: Gaussian's mean
and mode are the same)

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

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 LGM 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})

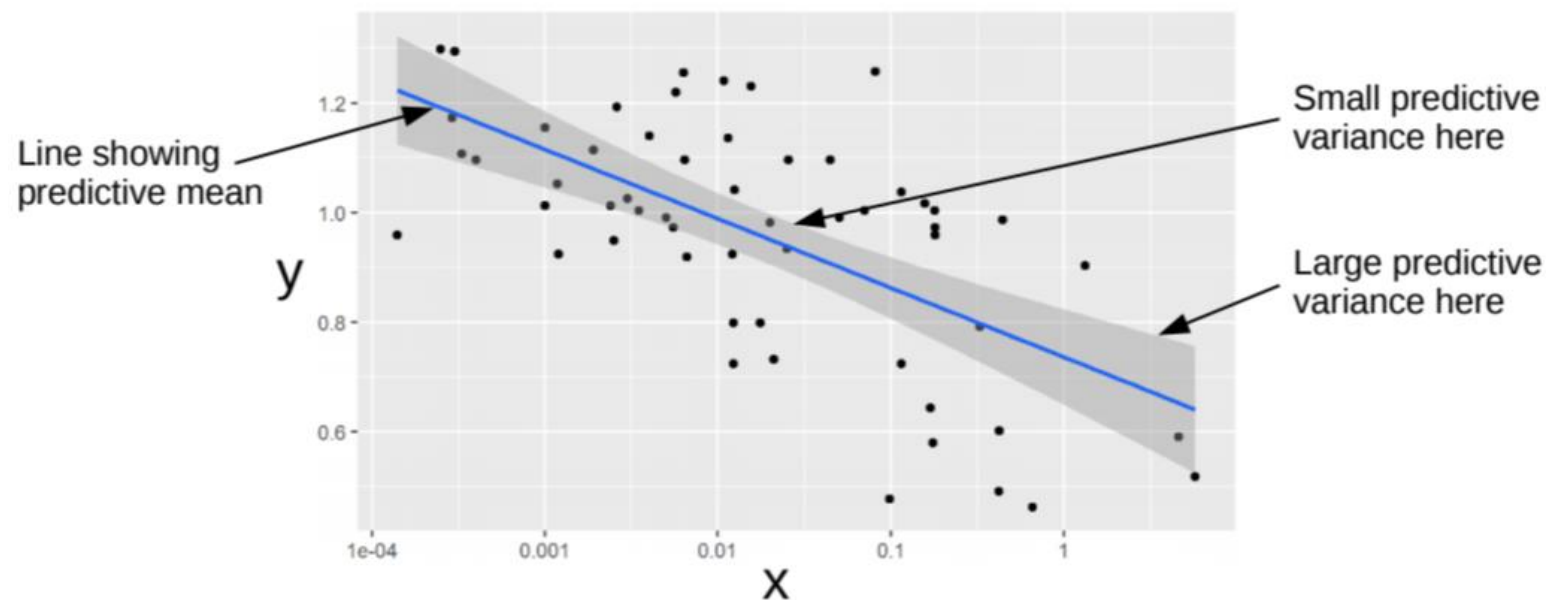
$$\begin{aligned} p(y_* | \mathbf{x}_*, \mathbf{w}_{MLE}) &= \mathcal{N}(\mathbf{w}_{MLE}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MLE prediction} \\ p(y_* | \mathbf{x}_*, \mathbf{w}_{MAP}) &= \mathcal{N}(\mathbf{w}_{MAP}^\top \mathbf{x}_*, \beta^{-1}) && \text{- MAP prediction} \end{aligned}$$

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

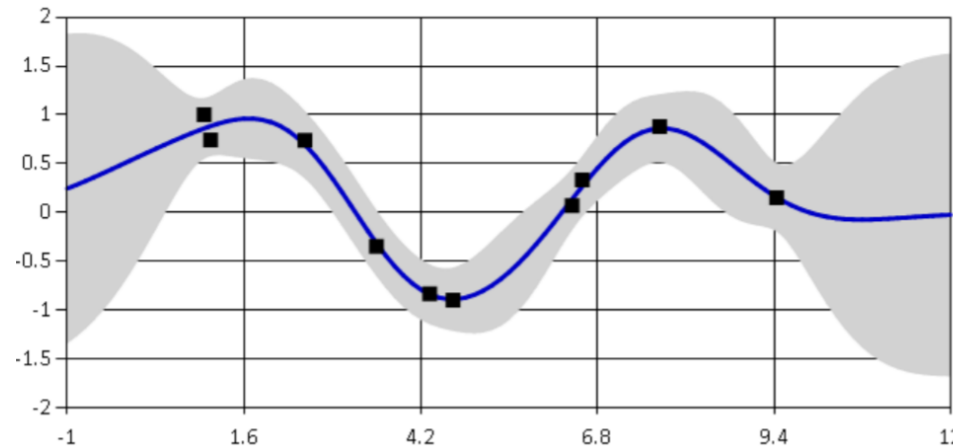
- Black dots are training examples



- 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 $\phi(\mathbf{x})$

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

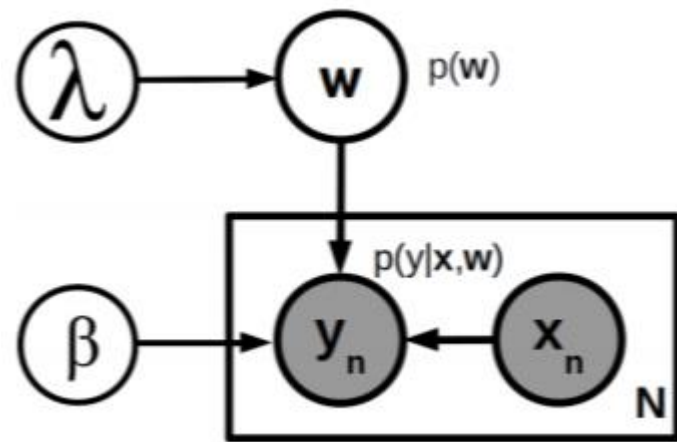
Can be pre-defined (e.g., replace a scalar x by polynomial mapping $[1, x, x^2]$) or extracted by a pretrained deep neural net

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



Estimating Hyperparameters via MLE-II

- 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

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

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

- Posterior and PPD computation is intractable.
- If we just want point estimates for (β, λ) then MLE-II is an option

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

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

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

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

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

Prob. Linear Regression: Some Other Variations

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

- Laplace distribution for the likelihood

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

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

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

Can even assume β_n to depend on input \mathbf{x}_n

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

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

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

This has the effect of having feature-specific regularization

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 $\mathcal{N}(\mathbf{w} | \mathbf{0}, \lambda^{-1} \mathbf{I})$ prior with spherical covariance

Diagonal precision/covariance matrix with λ_d 's along the columns of Λ

