

Probabilistic Models for Regression

Piyush Rai

Probabilistic Machine Learning (CS772A)

Aug 5, 2017

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right)$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \tau \mathbf{w}^\top \mathbf{w}$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \tau \mathbf{w}^\top \mathbf{w} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \tau \mathbf{I} \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right)$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \tau \mathbf{w}^\top \mathbf{w} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \tau \mathbf{I} \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

Background: Linear Regression

- Supervised Learning. Training data as N i.i.d. examples $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\mathbf{x}_n \in \mathbb{R}^D$, $y_n \in \mathbb{R}$
- Assume input-output relationship to approximately follow a **linear model** with parameters $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- $\mathbf{w} \in \mathbb{R}^D$ is also called the regression **weight vector** (the unknown we want to estimate)
 - Can think of the d^{th} component of \mathbf{w} , i.e. w_d , as **weight/importance** of d^{th} feature of the input
- We usually learn \mathbf{w} by optimizing a loss function or **regularized** loss function

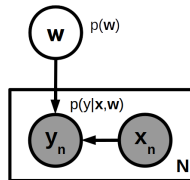
$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \tau \mathbf{w}^\top \mathbf{w} = \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top + \tau \mathbf{I} \right)^{-1} \left(\sum_{n=1}^N \mathbf{x}_n y_n \right) = (\mathbf{X}^\top \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

- Nice, convex objective with **unique solution**

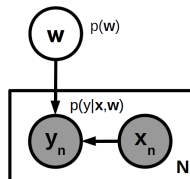
Probabilistic Linear Regression

- We would like to set up linear regression as a probabilistic model



Probabilistic Linear Regression

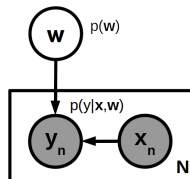
- We would like to set up linear regression as a probabilistic model



- We would specify the appropriate likelihood and prior distributions for this model

Probabilistic Linear Regression

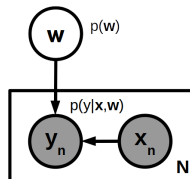
- We would like to set up linear regression as a probabilistic model



- We would specify the appropriate likelihood and prior distributions for this model
- We will then estimate the model parameters using MLE/MAP and fully Bayesian inference

Probabilistic Linear Regression

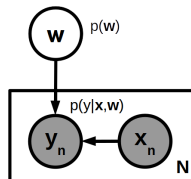
- We would like to set up linear regression as a probabilistic model



- We would specify the appropriate likelihood and prior distributions for this model
- We will then estimate the model parameters using MLE/MAP and fully Bayesian inference
- We will only be modeling the outputs y ; the inputs x will not be modeled (assumed fixed)

Probabilistic Linear Regression

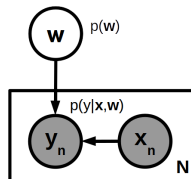
- We would like to set up linear regression as a probabilistic model



- We would specify the appropriate likelihood and prior distributions for this model
- We will then estimate the model parameters using MLE/MAP and fully Bayesian inference
- We will only be modeling the outputs y ; the inputs x will not be modeled (assumed fixed)
 - Note: In some supervised learning problems (especially classification), inputs are also modeled (more on this when we discuss “generative classification”)

Probabilistic Linear Regression

- We would like to set up linear regression as a probabilistic model

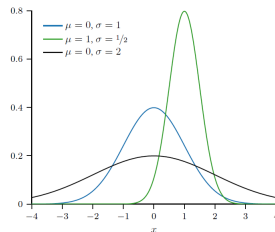


- We would specify the appropriate likelihood and prior distributions for this model
- We will then estimate the model parameters using MLE/MAP and fully Bayesian inference
- We will only be modeling the outputs y ; the inputs x will not be modeled (assumed fixed)
 - Note: In some supervised learning problems (especially classification), inputs are also modeled (more on this when we discuss “generative classification”)
- Only w and y_n ($n = 1, \dots, N$) are the random variables in this model (not x_n 's, which are fixed)

Background: Univariate Gaussian Distribution

- Distribution over real-valued scalar r.v. x
- Defined by a scalar **mean** μ and a scalar **variance** $\sigma^2 > 0$
- Mean: $\mathbb{E}[x] = \mu$, Variance: $\text{var}[x] = \sigma^2$
- Precision (inverse variance) $\beta = 1/\sigma^2$
- Distribution defined as

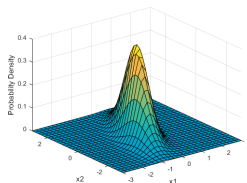
$$\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(x-\mu)^2}$$



Background: Multivariate Gaussian Distribution

- Distribution over a multivariate r.v. vector $\mathbf{x} \in \mathbb{R}^D$ of real numbers
- Defined by a **mean vector** $\boldsymbol{\mu} \in \mathbb{R}^D$ and a $D \times D$ **cov. matrix** $\boldsymbol{\Sigma}$ (its inverse is precision matrix)

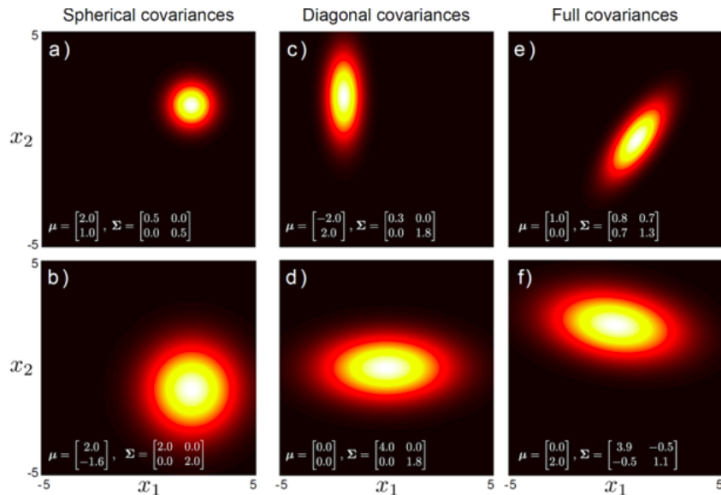
$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$



- The covariance matrix $\boldsymbol{\Sigma}$ must be symmetric and positive definite
 - All eigenvalues are positive
 - $\mathbf{z}^\top \boldsymbol{\Sigma} \mathbf{z} > 0$ for any real vector \mathbf{z}
- Inverse of the covariance matrix is known as the **precision matrix** $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}^{-1}$ (also positive definite)

Background: Multivariate Gaussian Distribution

The covariance matrix can be spherical, diagonal, or full



Picture courtesy: Computer vision: models, learning and inference (Simon Price)

Probabilistic Linear Regression: Likelihood Model

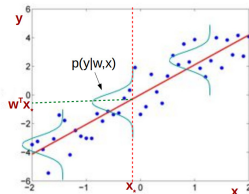
- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise

Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$

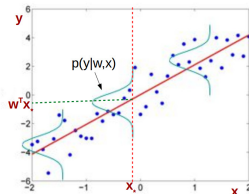
Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$



Probabilistic Linear Regression: Likelihood Model

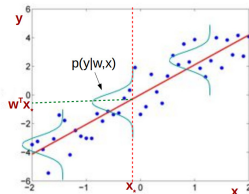
- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$



- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case

Probabilistic Linear Regression: Likelihood Model

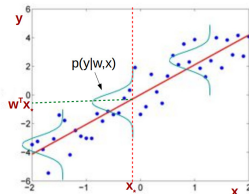
- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$



- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case
- Note: We are only modeling the outputs y ; the inputs \mathbf{x} are not being modeled (fixed)

Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$

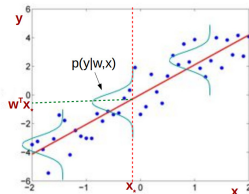


- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case
- Note: We are only modeling the outputs y ; the inputs \mathbf{x} are not being modeled (fixed)
- Denoting $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$ ($N \times D$), $\mathbf{y} = [y_1 \dots y_N]^\top$ ($N \times 1$), the overall likelihood for i.i.d. data

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w})$$

Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$

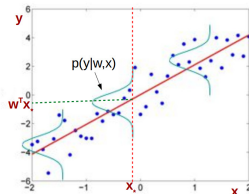


- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case
- Note: We are only modeling the outputs y ; the inputs \mathbf{x} are not being modeled (fixed)
- Denoting $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$ ($N \times D$), $\mathbf{y} = [y_1 \dots y_N]^\top$ ($N \times 1$), the overall likelihood for i.i.d. data

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$

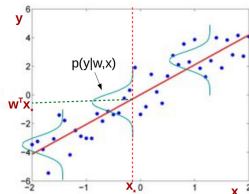


- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case
- Note: We are only modeling the outputs y ; the inputs \mathbf{x} are not being modeled (fixed)
- Denoting $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$ ($N \times D$), $\mathbf{y} = [y_1 \dots y_N]^\top$ ($N \times 1$), the overall likelihood for i.i.d. data

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2} \propto e^{-\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

Probabilistic Linear Regression: Likelihood Model

- Assume a model of the form $y_n = \mathbf{w}^\top \mathbf{x}_n + \epsilon_n$, where $\epsilon_n \sim \mathcal{N}(0, \beta^{-1})$ is Gaussian noise
- Thus each response y_n also has a Gaussian distribution: $p(y_n | \mathbf{x}_n, \mathbf{w}) = \mathcal{N}(\mathbf{w}^\top \mathbf{x}_n, \beta^{-1})$



- $p(y_n | \mathbf{x}_n, \mathbf{w})$ is our likelihood model - Gaussian distribution in this case
- Note: We are only modeling the outputs y ; the inputs \mathbf{x} are not being modeled (fixed)
- Denoting $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^\top$ ($N \times D$), $\mathbf{y} = [y_1 \dots y_N]^\top$ ($N \times 1$), the overall likelihood for i.i.d. data

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N p(y_n | \mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2} \propto e^{-\frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2} \propto e^{-\frac{\beta}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})}$$

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

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

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

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

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2}$$

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2}$$

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

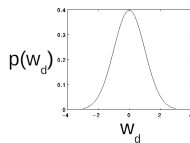
where λ denotes the precision of the Gaussian (how shrunk w_d is towards the mean of the Gaussian)

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

where λ denotes the precision of the Gaussian (how shrunk w_d is towards the mean of the Gaussian)

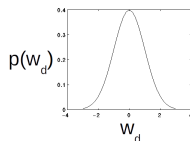


Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

where λ denotes the precision of the Gaussian (how shrunk w_d is towards the mean of the Gaussian)



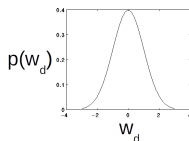
- The prior says that *a priori* we want the individual weights w_d to be small (close to zero)

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

where λ denotes the precision of the Gaussian (how shrunk w_d is towards the mean of the Gaussian)



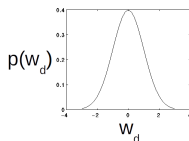
- The prior says that *a priori* we want the individual weights w_d to be small (close to zero)
 - Akin to imposing a regularizer on \mathbf{w} that promotes its elements to take small values

Probabilistic Linear Regression: Prior

- Let's assume a zero mean multivariate Gaussian prior on the weight vector $\mathbf{w} \in \mathbb{R}^D$

$$P(\mathbf{w}|\lambda) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D) = \prod_{d=1}^D \mathcal{N}(w_d|0, \lambda^{-1}) = \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} \propto e^{-\frac{\lambda}{2} \sum_{d=1}^D w_d^2} \propto e^{-\frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}}$$

where λ denotes the precision of the Gaussian (how shrunk w_d is towards the mean of the Gaussian)



- The prior says that *a priori* we want the individual weights w_d to be small (close to zero)
 - Akin to imposing a regularizer on \mathbf{w} that promotes its elements to take small values
 - We'll see that the Gaussian prior exactly corresponds to having an ℓ_2 (squared Euclidean norm) regularizer on \mathbf{w} , of the form $\mathbf{w}^\top \mathbf{w} = \sum_{d=1}^D w_d^2$

MLE for Probabilistic Linear Regression

- The negative log-likelihood will be

$$NLL(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\log \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

MLE for Probabilistic Linear Regression

- The negative log-likelihood will be

$$NLL(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\log \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

- Simplifying further and ignoring the constants w.r.t. \mathbf{w} , we get

$$NLL(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

MLE for Probabilistic Linear Regression

- The negative log-likelihood will be

$$NLL(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\log \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

- Simplifying further and ignoring the constants w.r.t. \mathbf{w} , we get

$$NLL(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Note that this objective has the same form as unregularized linear regression

MLE for Probabilistic Linear Regression

- The negative log-likelihood will be

$$NLL(\mathbf{w}) = -\log p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = -\log \prod_{n=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2} (y_n - \mathbf{w}^\top \mathbf{x}_n)^2}$$

- Simplifying further and ignoring the constants w.r.t. \mathbf{w} , we get

$$NLL(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Note that this objective has the same form as unregularized linear regression
- Therefore the MLE solution will be $\hat{\mathbf{w}}_{MLE} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

MAP Estimation for Probabilistic Linear Regression

- The negative log-prior will be (ignoring constants w.r.t. \mathbf{w})

$$-\log p(\mathbf{w}|\lambda) = -\log \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} = \frac{\lambda}{2} \sum_{d=1}^D w_d^2 = \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

MAP Estimation for Probabilistic Linear Regression

- The negative log-prior will be (ignoring constants w.r.t. \mathbf{w})

$$-\log p(\mathbf{w}|\lambda) = -\log \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} = \frac{\lambda}{2} \sum_{d=1}^D w_d^2 = \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Therefore the MAP solution $\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) - \log p(\mathbf{w}|\lambda)$ will be

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \sum_{d=1}^D w_d^2$$

MAP Estimation for Probabilistic Linear Regression

- The negative log-prior will be (ignoring constants w.r.t. \mathbf{w})

$$-\log p(\mathbf{w}|\lambda) = -\log \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} = \frac{\lambda}{2} \sum_{d=1}^D w_d^2 = \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Therefore the MAP solution $\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) - \log p(\mathbf{w}|\lambda)$ will be

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \sum_{d=1}^D w_d^2$$

- Can also write the above in the following form

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{\beta} \mathbf{w}^\top \mathbf{w}$$

MAP Estimation for Probabilistic Linear Regression

- The negative log-prior will be (ignoring constants w.r.t. \mathbf{w})

$$-\log p(\mathbf{w}|\lambda) = -\log \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} = \frac{\lambda}{2} \sum_{d=1}^D w_d^2 = \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Therefore the MAP solution $\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) - \log p(\mathbf{w}|\lambda)$ will be

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \sum_{d=1}^D w_d^2$$

- Can also write the above in the following form

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{\beta} \mathbf{w}^\top \mathbf{w}$$

- Note that this objective has the same form as regularized (“ridge”) linear regression with $\tau = \frac{\lambda}{\beta}$

MAP Estimation for Probabilistic Linear Regression

- The negative log-prior will be (ignoring constants w.r.t. \mathbf{w})

$$-\log p(\mathbf{w}|\lambda) = -\log \prod_{d=1}^D \sqrt{\frac{\lambda}{2\pi}} e^{-\frac{\lambda}{2} w_d^2} = \frac{\lambda}{2} \sum_{d=1}^D w_d^2 = \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w}$$

- Therefore the MAP solution $\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) - \log p(\mathbf{w}|\lambda)$ will be

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \frac{\beta}{2} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{2} \sum_{d=1}^D w_d^2$$

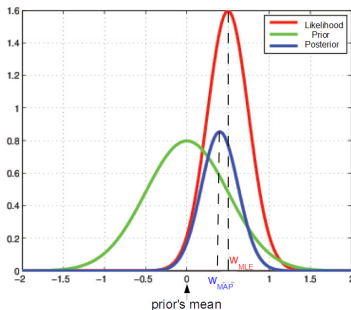
- Can also write the above in the following form

$$\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \frac{\lambda}{\beta} \mathbf{w}^\top \mathbf{w}$$

- Note that this objective has the same form as regularized (“ridge”) linear regression with $\tau = \frac{\lambda}{\beta}$
- Therefore the MAP solution will be $\hat{\mathbf{w}}_{MAP} = \arg \min_{\mathbf{w}} NLL(\mathbf{w}) = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

MLE vs MAP: An Illustration

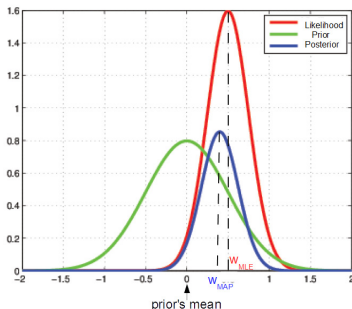
- \mathbf{w}_{MAP} is a compromise between prior's mean (zero in this case) and \mathbf{w}_{MLE}



- Doing MAP shrinks the estimate of \mathbf{w} towards the prior's mean

MLE vs MAP: An Illustration

- \mathbf{w}_{MAP} is a compromise between prior's mean (zero in this case) and \mathbf{w}_{MLE}



- Doing MAP shrinks the estimate of \mathbf{w} towards the prior's mean
- Note that a “peaked” prior (large precision) will pull the solution more towards the prior mean

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}}_{MAP} = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}}_{MAP} = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\hat{\mathbf{w}}_{MLE} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\hat{\mathbf{w}}_{MAP} = (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)
- MAP uses a Gaussian prior on \mathbf{w} (which acts as a regularizer on \mathbf{w})

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\begin{aligned}\hat{\mathbf{w}}_{MLE} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \hat{\mathbf{w}}_{MAP} &= (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)
- MAP uses a Gaussian prior on \mathbf{w} (which acts as a regularizer on \mathbf{w})
- As the prior's precision goes to zero, $\hat{\mathbf{w}}_{MAP}$ approaches $\hat{\mathbf{w}}_{MLE}$

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\begin{aligned}\hat{\mathbf{w}}_{MLE} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \hat{\mathbf{w}}_{MAP} &= (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)
- MAP uses a Gaussian prior on \mathbf{w} (which acts as a regularizer on \mathbf{w})
- As the prior's precision goes to zero, $\hat{\mathbf{w}}_{MAP}$ approaches $\hat{\mathbf{w}}_{MLE}$
 - Likewise, the higher the prior's precision λ is, the more *regularized* the solution will be

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\begin{aligned}\hat{\mathbf{w}}_{MLE} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \hat{\mathbf{w}}_{MAP} &= (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)
- MAP uses a Gaussian prior on \mathbf{w} (which acts as a regularizer on \mathbf{w})
- As the prior's precision goes to zero, $\hat{\mathbf{w}}_{MAP}$ approaches $\hat{\mathbf{w}}_{MLE}$
 - Likewise, the higher the prior's precision λ is, the more *regularized* the solution will be
- Prediction for a new test input \mathbf{x}_* will be

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

MLE vs MAP for Probabilistic Linear Regression

- The MLE and MAP solutions are

$$\begin{aligned}\hat{\mathbf{w}}_{MLE} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \\ \hat{\mathbf{w}}_{MAP} &= (\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

- MLE doesn't use a prior on \mathbf{w} (akin to unregularized linear regression)
- MAP uses a Gaussian prior on \mathbf{w} (which acts as a regularizer on \mathbf{w})
- As the prior's precision goes to zero, $\hat{\mathbf{w}}_{MAP}$ approaches $\hat{\mathbf{w}}_{MLE}$
 - Likewise, the higher the prior's precision λ is, the more *regularized* the solution will be
- Prediction for a new test input \mathbf{x}_* will be

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

- Predictive means are different but predictive variances are the same (and constant for all $\mathbf{x}!$)

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)
- The posterior over the weight vector \mathbf{w} (with β and λ known)

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

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)
- The posterior over the weight vector \mathbf{w} (with β and λ known)

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

- Computing $P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$ (like Bernoulli-Beta case, doing it only upto proportionality constant)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto P(\mathbf{w}|\lambda)P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$$

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)
- The posterior over the weight vector \mathbf{w} (with β and λ known)

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

- Computing $P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$ (like Bernoulli-Beta case, doing it only upto proportionality constant)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto P(\mathbf{w}|\lambda)P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$$

- After some algebra, this gets simplified into the following (proof on the next two slides)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (\text{The posterior must be Gaussian due to conjugacy})$$

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)
- The posterior over the weight vector \mathbf{w} (with β and λ known)

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

- Computing $P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$ (like Bernoulli-Beta case, doing it only upto proportionality constant)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto P(\mathbf{w}|\lambda)P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$$

- After some algebra, this gets simplified into the following (proof on the next two slides)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (\text{The posterior must be Gaussian due to conjugacy})$$

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

Inferring the Posterior Distribution (fully Bayesian Inference)

- Inferring the full posterior is straightforward if the hyperparams β and λ to be known/fixed
 - Basically, the conjugacy helps here (Gaussian prior is **conjugate** to Gaussian likelihood)
- The posterior over the weight vector \mathbf{w} (with β and λ known)

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

- Computing $P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda)$ (like Bernoulli-Beta case, doing it only upto proportionality constant)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto P(\mathbf{w}|\lambda)P(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$$

- After some algebra, this gets simplified into the following (proof on the next two slides)

$$P(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad (\text{The posterior must be Gaussian due to conjugacy})$$

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

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\beta \sum_{n=1}^N y_n \mathbf{x}_n \right) = \boldsymbol{\Sigma} (\beta \mathbf{X}^{\top} \mathbf{y}) = (\mathbf{X}^{\top} \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

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

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) &\propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) &= \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \\ &\propto \exp\left(-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right)\exp\left(-\frac{\beta}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^\top(\mathbf{y}-\mathbf{X}\mathbf{w})\right) \end{aligned}$$

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) &\propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) &= \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \\ &\propto \exp\left(-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right)\exp\left(-\frac{\beta}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^\top(\mathbf{y}-\mathbf{X}\mathbf{w})\right) \\ &= \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{y}^\top\mathbf{y} + \mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \end{aligned}$$

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) &\propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \\ &\propto \exp\left(-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right)\exp\left(-\frac{\beta}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^\top(\mathbf{y}-\mathbf{X}\mathbf{w})\right) \\ &= \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{y}^\top\mathbf{y} + \mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \\ &\propto \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \end{aligned}$$

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) &\propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \\ &\propto \exp\left(-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right)\exp\left(-\frac{\beta}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^\top(\mathbf{y}-\mathbf{X}\mathbf{w})\right) \\ &= \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{y}^\top\mathbf{y} + \mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \\ &\propto \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \\ &= \exp\left[-\frac{1}{2}\left(\mathbf{w}^\top(\lambda\mathbf{I}_D + \beta\mathbf{X}^\top\mathbf{X})\mathbf{w} - 2\beta\mathbf{w}^\top\mathbf{X}^\top\mathbf{y}\right)\right] \end{aligned}$$

The “Completing The Square” Trick for Gaussian Posterior

- Plugging in the respective distributions for $p(\mathbf{w}|\lambda)$ and $p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)$, we will get

$$\begin{aligned} p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) &\propto p(\mathbf{w}|\lambda)p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \lambda^{-1}\mathbf{I}_D)\mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \beta^{-1}\mathbf{I}_N) \\ &\propto \exp\left(-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w}\right)\exp\left(-\frac{\beta}{2}(\mathbf{y}-\mathbf{X}\mathbf{w})^\top(\mathbf{y}-\mathbf{X}\mathbf{w})\right) \\ &= \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{y}^\top\mathbf{y} + \mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \\ &\propto \exp\left[-\frac{\lambda}{2}\mathbf{w}^\top\mathbf{w} - \frac{\beta}{2}(\mathbf{w}^\top\mathbf{X}^\top\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top\mathbf{X}^\top\mathbf{y})\right] \\ &= \exp\left[-\frac{1}{2}\left(\mathbf{w}^\top(\lambda\mathbf{I}_D + \beta\mathbf{X}^\top\mathbf{X})\mathbf{w} - 2\beta\mathbf{w}^\top\mathbf{X}^\top\mathbf{y}\right)\right] \end{aligned}$$

- We will now try to bring the exponent into a quadratic form to see if it corresponds to some Gaussian. So basically, we will use the “complete the square” trick

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$
- Let's see if we can bring the above posterior into the form of the following Gaussian

$$\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left[-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right] = \exp \left[-\frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right]$$

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$
- Let's see if we can bring the above posterior into the form of the following Gaussian

$$\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left[-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right] = \exp \left[-\frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right]$$

- Let's multiply and divide $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$ by $\exp \left[-\frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right]$

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$
- Let's see if we can bring the above posterior into the form of the following Gaussian

$$\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left[-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right] = \exp \left[-\frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right]$$

- Let's multiply and divide $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$ by $\exp \left[-\frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right]$
- This gives the following up to a prop. constant (remember $\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ is constant w.r.t. \mathbf{w}):

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \right]$$

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$
- Let's see if we can bring the above posterior into the form of the following Gaussian

$$\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left[-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right] = \exp \left[-\frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right]$$

- Let's multiply and divide $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$ by $\exp \left[-\frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right]$
- This gives the following up to a prop. constant (remember $\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ is constant w.r.t. \mathbf{w}):

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \right]$$

- Finally comparing with the expression of $\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ we can see that

$$\begin{aligned} \boldsymbol{\Sigma} &= (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X})^{-1} \\ \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} &= \beta \mathbf{X}^\top \mathbf{y} \Rightarrow \boldsymbol{\mu} = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) \end{aligned}$$

The “Completing The Square” Trick for Gaussian Posterior

- So we had.. $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$
- Let's see if we can bring the above posterior into the form of the following Gaussian

$$\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \exp \left[-\frac{1}{2} (\mathbf{w} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{w} - \boldsymbol{\mu}) \right] = \exp \left[-\frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \mathbf{w} - 2\mathbf{w}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}) \right]$$

- Let's multiply and divide $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} \right) \right]$ by $\exp \left[-\frac{1}{2} \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right]$
- This gives the following up to a prop. constant (remember $\boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}$ is constant w.r.t. \mathbf{w}):

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) \propto \exp \left[-\frac{1}{2} \left(\mathbf{w}^\top (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X}) \mathbf{w} - 2\beta \mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\mu}^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} \right) \right]$$

- Finally comparing with the expression of $\mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ we can see that

$$\begin{aligned} \boldsymbol{\Sigma} &= (\lambda \mathbf{I}_D + \beta \mathbf{X}^\top \mathbf{X})^{-1} \\ \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} &= \beta \mathbf{X}^\top \mathbf{y} \Rightarrow \boldsymbol{\mu} = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) \end{aligned}$$

- Note: The above expression for the posterior can also be directly obtained using properties of Gaussian distributions (we will see those in the coming lectures)

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\beta \sum_{n=1}^N y_n \mathbf{x}_n \right) = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\beta \sum_{n=1}^N y_n \mathbf{x}_n \right) = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Note the form of the mean $\boldsymbol{\mu}$ of the posterior. Same as the MAP estimate \mathbf{w}_{MAP}

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\beta \sum_{n=1}^N y_n \mathbf{x}_n \right) = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Note the form of the mean $\boldsymbol{\mu}$ of the posterior. Same as the MAP estimate \mathbf{w}_{MAP}
 - Reason: For a Gaussian distribution, mean is the same as the mode (maxima of the distribution)

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left(\beta \sum_{n=1}^N y_n \mathbf{x}_n \right) = \boldsymbol{\Sigma} (\beta \mathbf{X}^\top \mathbf{y}) = \left(\mathbf{X}^\top \mathbf{X} + \frac{\lambda}{\beta} \mathbf{I}_D \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- Note the form of the mean $\boldsymbol{\mu}$ of the posterior. Same as the MAP estimate \mathbf{w}_{MAP}
 - Reason: For a Gaussian distribution, mean is the same as the mode (maxima of the distribution)
- Setting λ to 0 makes $\boldsymbol{\mu}$ equivalent to the MLE solution

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

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

- Note the form of the mean $\boldsymbol{\mu}$ of the posterior. Same as the MAP estimate \mathbf{w}_{MAP}
 - Reason: For a Gaussian distribution, mean is the same as the mode (maxima of the distribution)
- Setting λ to 0 makes $\boldsymbol{\mu}$ equivalent to the MLE solution
- However, MLE and MAP are only point estimates whereas we now have the full posterior over \mathbf{w}

Information within the Posterior

- As we saw, the posterior over the weight vector \mathbf{w}

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

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

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

- Note the form of the mean $\boldsymbol{\mu}$ of the posterior. Same as the MAP estimate \mathbf{w}_{MAP}
 - Reason: For a Gaussian distribution, mean is the same as the mode (maxima of the distribution)
- Setting λ to 0 makes $\boldsymbol{\mu}$ equivalent to the MLE solution
- However, MLE and MAP are only point estimates whereas we now have the full posterior over \mathbf{w}
- Therefore, the solution given by the full posterior in a way subsumes MAP/MLE solutions

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)
- The [posterior predictive distribution](#) in this case will be

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

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)
- The posterior predictive distribution in this case will be

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

- Both terms on RHS are Gaussians: $p(y_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(y_*|\mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$ and $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)
- The posterior predictive distribution in this case will be

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

- Both terms on RHS are Gaussians: $p(y_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(y_*|\mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$ and $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- The result and therefore the posterior predictive will be another Gaussian (we will see proof later)

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

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)
- The posterior predictive distribution in this case will be

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

- Both terms on RHS are Gaussians: $p(y_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(y_*|\mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$ and $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- The result and therefore the posterior predictive will be another Gaussian (we will see proof later)

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

- In contrast, MLE and MAP can only give us “plug-in” predictions (using the point estimate of \mathbf{w})

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

Posterior Predictive Distribution

- Now we want to use the posterior over \mathbf{w} to make predictions (response y_* for a new input \mathbf{x}_*)
- The [posterior predictive distribution](#) in this case will be

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

- Both terms on RHS are Gaussians: $p(y_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(y_*|\mathbf{w}^\top \mathbf{x}_*, \beta^{-1})$ and $p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \beta, \lambda) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- The result and therefore the posterior predictive will be another Gaussian (we will see proof later)

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

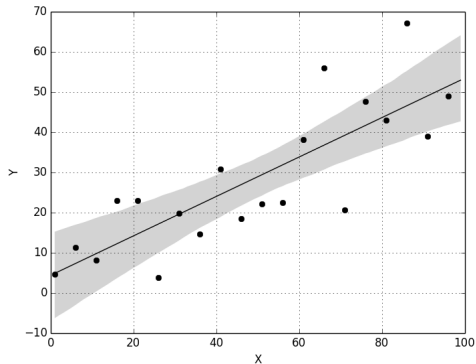
- In contrast, MLE and MAP can only give us “plug-in” predictions (using the point estimate of \mathbf{w})

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

- 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



Regions with more training examples have smaller predictive variance

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation
- However, we can also learn these hyperparameters

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation
- However, we can also learn these hyperparameters
 - Marginal likelihood maximization (also called MLE-II) is one such approach

$$\hat{\lambda}, \hat{\beta} = \arg \max_{\lambda, \beta} p(\mathbf{y} | \lambda, \beta, \mathbf{X})$$

where $p(\mathbf{y} | \lambda, \beta, \mathbf{X}) = \int_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \lambda, \beta, \mathbf{X}) p(\mathbf{w} | \lambda) d\mathbf{w}$ is the marginal likelihood

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation
- However, we can also learn these hyperparameters
 - Marginal likelihood maximization (also called MLE-II) is one such approach

$$\hat{\lambda}, \hat{\beta} = \arg \max_{\lambda, \beta} p(\mathbf{y} | \lambda, \beta, \mathbf{X})$$

where $p(\mathbf{y} | \lambda, \beta, \mathbf{X}) = \int_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \lambda, \beta, \mathbf{X}) p(\mathbf{w} | \lambda) d\mathbf{w}$ is the marginal likelihood

- Note: This optimization problem may not have a closed-form solution (iterative methods can be used)

Hyperparameter Estimation

- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation
- However, we can also learn these hyperparameters
 - Marginal likelihood maximization (also called MLE-II) is one such approach

$$\hat{\lambda}, \hat{\beta} = \arg \max_{\lambda, \beta} p(\mathbf{y} | \lambda, \beta, \mathbf{X})$$

where $p(\mathbf{y} | \lambda, \beta, \mathbf{X}) = \int_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \lambda, \beta, \mathbf{X}) p(\mathbf{w} | \lambda) d\mathbf{w}$ is the marginal likelihood

- Note: This optimization problem may not have a closed-form solution (iterative methods can be used)
- Marginal likelihood is also known as “evidence”

Hyperparameter Estimation

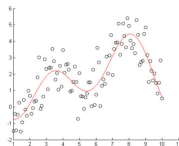
- So far we assumed that only \mathbf{w} is the unknown and the hyperparameters λ and β are known
- One way to choose the “optimal” values of the hyperparams is cross-validation
- However, we can also learn these hyperparameters
 - Marginal likelihood maximization (also called MLE-II) is one such approach

$$\hat{\lambda}, \hat{\beta} = \arg \max_{\lambda, \beta} p(\mathbf{y} | \lambda, \beta, \mathbf{X})$$

where $p(\mathbf{y} | \lambda, \beta, \mathbf{X}) = \int_{\mathbf{w}} p(\mathbf{y} | \mathbf{w}, \lambda, \beta, \mathbf{X}) p(\mathbf{w} | \lambda) d\mathbf{w}$ is the marginal likelihood

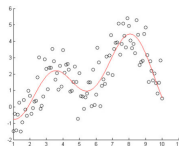
- Note: This optimization problem may not have a closed-form solution (iterative methods can be used)
- Marginal likelihood is also known as “evidence”
- This procedure of hyperparameter estimation is also called “evidence maximization”

Probabilistic (Nonlinear) Regression



- Many ways to extend the probabilistic linear regression model to nonlinear regression.

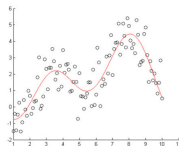
Probabilistic (Nonlinear) Regression



- Many ways to extend the probabilistic linear regression model to nonlinear regression. E.g.,
 - Replace \mathbf{x} by a **pre-defined higher-dim. feature mapping** $\phi(\mathbf{x})$ (e.g. $\mathbf{x} \rightarrow (1, \mathbf{x}, \mathbf{x}^2)$) and apply the linear model on the new feature mapped version of the inputs

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

Probabilistic (Nonlinear) Regression



- Many ways to extend the probabilistic linear regression model to nonlinear regression. E.g.,
 - Replace \mathbf{x} by a **pre-defined higher-dim. feature mapping** $\phi(\mathbf{x})$ (e.g. $\mathbf{x} \rightarrow (1, \mathbf{x}, \mathbf{x}^2)$) and apply the linear model on the new feature mapped version of the inputs

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

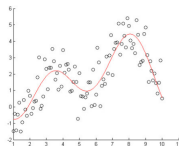
.. also related to kernel methods

- Replace the linear model $\mathbf{w}^\top \mathbf{x}$ by a **nonlinear function** f (e.g., a neural network)

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

.. more on this when we discuss deep probabilistic models

Probabilistic (Nonlinear) Regression



- Many ways to extend the probabilistic linear regression model to nonlinear regression. E.g.,
 - Replace \mathbf{x} by a **pre-defined higher-dim. feature mapping** $\phi(\mathbf{x})$ (e.g. $\mathbf{x} \rightarrow (1, \mathbf{x}, \mathbf{x}^2)$) and apply the linear model on the new feature mapped version of the inputs

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

.. also related to kernel methods

- Replace the linear model $\mathbf{w}^\top \mathbf{x}$ by a **nonlinear function** f (e.g., a neural network)

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

.. more on this when we discuss deep probabilistic models

- Use **Gaussian Processes**