



UNIVERSIDADE
FEDERAL DO CEARÁ



LogIA

Aprendizagem de Máquina Probabilística

César Lincoln Cavalcante Mattos

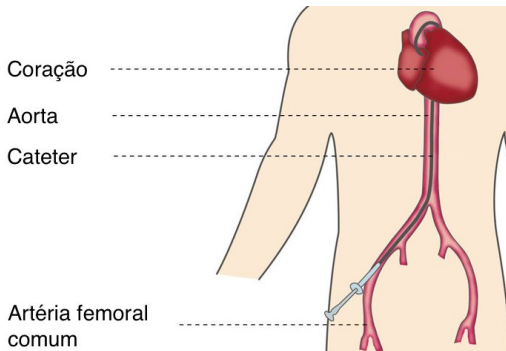
2025

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

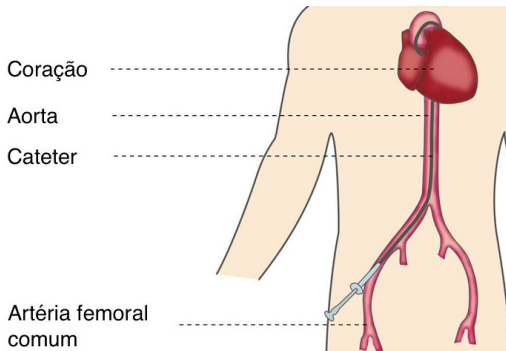
Regressão linear simples

- **Cateterismo cardíaco:** procedimento de inserção de um cateter (usualmente pela artéria femoral) para diagnosticar problemas de obstrução no coração.



Regressão linear simples

- **Cateterismo cardíaco:** procedimento de inserção de um cateter (usualmente pela artéria femoral) para diagnosticar problemas de obstrução no coração.



- **Problema:** Dada a **altura** de um paciente, qual o **comprimento** do cateter necessário para alcançar seu coração?

Regressão linear simples

- Considere a tabela a seguir relacionando alturas de jovens pacientes e comprimentos do cateter correspondente:

Altura (m)	Comprimento (cm)
1.087	37
1.613	50
0.953	34
1.003	36
1.156	43
0.978	28
1.092	37
0.572	20
0.940	34
0.597	30
0.838	38
1.473	47

Regressão linear simples

- Considere a tabela a seguir relacionando alturas de jovens pacientes e comprimentos do cateter correspondente:

Altura (m)	Comprimento (cm)
1.087	37
1.613	50
0.953	34
1.003	36
1.156	43
0.978	28
1.092	37
0.572	20
0.940	34
0.597	30
0.838	38
1.473	47

- **Problema:** Dado uma altura **não presente** na tabela, qual deverá ser o comprimento do cateter?

Regressão linear simples

- A coluna **Altura** é a **entrada** do nosso modelo.
- A coluna **Comprimento** é a **saída** do nosso modelo.
- Nosso **conjunto de dados** é formado por 12 alturas e 12 comprimentos correspondentes.

Regressão linear simples

- A coluna **Altura** é a **entrada** do nosso modelo.
- A coluna **Comprimento** é a **saída** do nosso modelo.
- Nosso **conjunto de dados** é formado por 12 alturas e 12 comprimentos correspondentes.
- Matematicamente, temos:

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_2, y_2)\} = \{(x_i, y_i)\}_{i=1}^{12},$$

em que x_i é a i -ésima entrada e y_i é a i -ésima saída.

Regressão linear simples

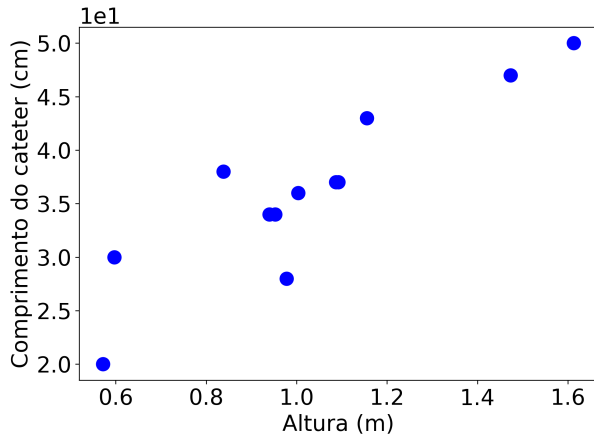
- A coluna **Altura** é a **entrada** do nosso modelo.
- A coluna **Comprimento** é a **saída** do nosso modelo.
- Nosso **conjunto de dados** é formado por 12 alturas e 12 comprimentos correspondentes.
- Matematicamente, temos:

$$\mathcal{D} = \{(x_1, y_1), \dots, (x_2, y_2)\} = \{(x_i, y_i)\}_{i=1}^{12},$$

em que x_i é a i -ésima entrada e y_i é a i -ésima saída.

- **Objetivo:** Encontrar uma relação entre x_i e y_i que forneça uma predição \hat{y}_i o mais próximo possível da saída real y_i .

Regressão linear simples



Regressão linear simples

Terminologia

- **Atributo (feature)**: Uma dada característica de um padrão.

Regressão linear simples

Terminologia

- **Atributo (feature)**: Uma dada característica de um padrão.
- **Padrão (pattern)**: Um vetor de atributos que representa um exemplo.

Regressão linear simples

Terminologia

- **Atributo (feature)**: Uma dada característica de um padrão.
- **Padrão (pattern)**: Um vetor de atributos que representa um exemplo.
- **Modelo**: Uma função que expressa a relação entre um padrão de entrada e sua saída correspondente.

Regressão linear simples

Terminologia

- **Atributo (feature)**: Uma dada característica de um padrão.
- **Padrão (pattern)**: Um vetor de atributos que representa um exemplo.
- **Modelo**: Uma função que expressa a relação entre um padrão de entrada e sua saída correspondente.
- **Função custo (ou função objetivo)**: Indica o quão mal (ou o quão bem) um modelo aproxima os dados disponíveis.

Regressão linear simples

Terminologia

- **Atributo (feature)**: Uma dada característica de um padrão.
- **Padrão (pattern)**: Um vetor de atributos que representa um exemplo.
- **Modelo**: Uma função que expressa a relação entre um padrão de entrada e sua saída correspondente.
- **Função custo (ou função objetivo)**: Indica o quão mal (ou o quão bem) um modelo aproxima os dados disponíveis.
- **Parâmetros**: Variáveis que caracterizam o modelo proposto.

Regressão linear simples

Terminologia

- **Atributo (feature):** Uma dada característica de um padrão.
- **Padrão (pattern):** Um vetor de atributos que representa um exemplo.
- **Modelo:** Uma função que expressa a relação entre um padrão de entrada e sua saída correspondente.
- **Função custo (ou função objetivo):** Indica o quão mal (ou o quão bem) um modelo aproxima os dados disponíveis.
- **Parâmetros:** Variáveis que caracterizam o modelo proposto.
- **Risco empírico:** Estimativa do risco (custo) obtida a partir dos dados disponíveis.

Regressão linear simples

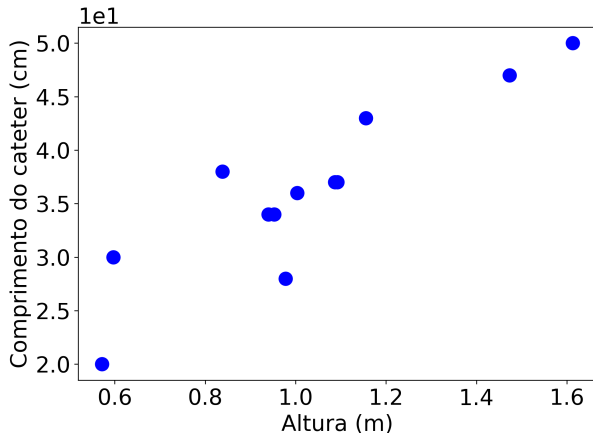
Terminologia

- **Atributo (feature):** Uma dada característica de um padrão.
- **Padrão (pattern):** Um vetor de atributos que representa um exemplo.
- **Modelo:** Uma função que expressa a relação entre um padrão de entrada e sua saída correspondente.
- **Função custo (ou função objetivo):** Indica o quão mal (ou o quão bem) um modelo aproxima os dados disponíveis.
- **Parâmetros:** Variáveis que caracterizam o modelo proposto.
- **Risco empírico:** Estimativa do risco (custo) obtida a partir dos dados disponíveis.
- **Otimização (ou treinamento, aprendizagem):** Algoritmo de obtenção dos parâmetros do modelo que minimizem uma função custo (ou maximizem uma função objetivo).

Regressão linear simples

- Considere uma relação linear entre x_i (entrada do modelo) e \hat{y}_i (saída do modelo):

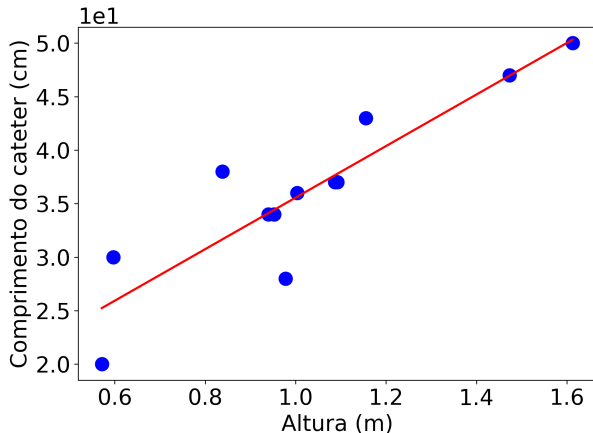
$$\hat{y}_i = w_0 + w_1 x_i.$$



Regressão linear simples

- Considere uma relação linear entre x_i (entrada do modelo) e \hat{y}_i (saída do modelo):

$$\hat{y}_i = w_0 + w_1 x_i.$$



Regressão linear simples

- Escolhemos uma função custo quadrática para os erros obtidos pelo modelo:

$$\mathcal{J}(w_0, w_1) = \frac{1}{2N} \sum_{i=1}^N e_i^2,$$
$$e_i = y_i - \hat{y}_i.$$

Regressão linear simples

- Escolhemos uma função custo quadrática para os erros obtidos pelo modelo:

$$\mathcal{J}(w_0, w_1) = \frac{1}{2N} \sum_{i=1}^N e_i^2,$$
$$e_i = y_i - \hat{y}_i.$$

- Esse custo é chamado **Erro Quadrático Médio (MSE, Mean Squared Error)**.

Regressão linear simples

- Desejamos minimizar a função custo em relação aos parâmetros do modelo:

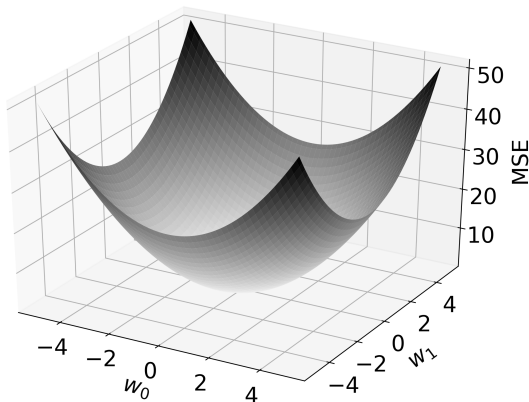
$$\min_{w_0, w_1} \mathcal{J}(w_0, w_1)$$

$$\min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

$$\min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

Regressão linear simples

$$\min_{w_0, w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$



Regressão linear simples

- Escolhemos valores iniciais para w_0 e w_1 .
- Movimentamos os parâmetros na direção que diminui a função custo $\mathcal{J}(w_0, w_1)$:

$$w_0 \leftarrow w_0 - \alpha \frac{\partial \mathcal{J}}{\partial w_0},$$
$$w_1 \leftarrow w_1 - \alpha \frac{\partial \mathcal{J}}{\partial w_1}$$

- $\alpha > 0$ é um **passo de aprendizado**.

Regressão linear simples

$$\frac{\partial \mathcal{J}}{\partial w_0} = \frac{\partial}{\partial w_0} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

$$\frac{\partial \mathcal{J}}{\partial w_0} = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)(-1)$$

$$\frac{\partial \mathcal{J}}{\partial w_0} = -\frac{1}{N} \sum_{i=1}^N e_i$$

Regressão linear simples

$$\frac{\partial \mathcal{J}}{\partial w_0} = \frac{\partial}{\partial w_0} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

$$\frac{\partial \mathcal{J}}{\partial w_0} = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)(-1)$$

$$\frac{\partial \mathcal{J}}{\partial w_0} = -\frac{1}{N} \sum_{i=1}^N e_i$$

- Logo:

$$w_0 \leftarrow w_0 + \alpha \frac{1}{N} \sum_{i=1}^N e_i$$

Regressão linear simples

$$\frac{\partial \mathcal{J}}{\partial w_1} = \frac{\partial}{\partial w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

$$\frac{\partial \mathcal{J}}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)(-x_i)$$

$$\frac{\partial \mathcal{J}}{\partial w_1} = -\frac{1}{N} \sum_{i=1}^N e_i x_i$$

Regressão linear simples

$$\frac{\partial \mathcal{J}}{\partial w_1} = \frac{\partial}{\partial w_1} \frac{1}{2N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)^2$$

$$\frac{\partial \mathcal{J}}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N (y_i - w_0 - w_1 x_i)(-x_i)$$

$$\frac{\partial \mathcal{J}}{\partial w_1} = -\frac{1}{N} \sum_{i=1}^N e_i x_i$$

- Logo:

$$w_1 \leftarrow w_1 + \alpha \frac{1}{N} \sum_{i=1}^N e_i x_i$$

Regressão linear simples

Gradiente Descendente (GD, gradient descent)

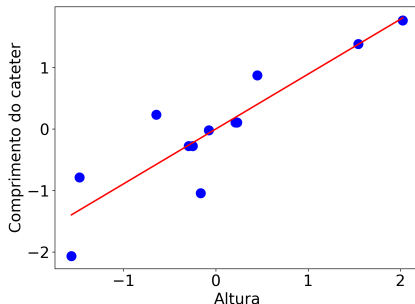
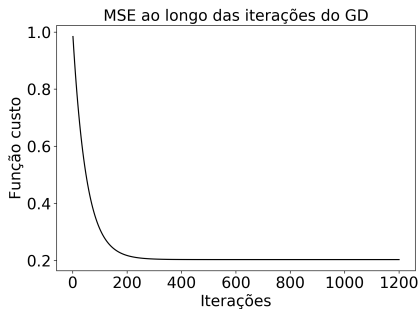
- 1 Escolha um valor α positivo e pequeno.
- 2 Inicialize os parâmetros do modelo na iteração $t = 0$.
- 3 Repita por diversas iterações (épocas):
 - 1 $t \leftarrow t + 1$;
 - 2 Calcule os erros do modelo:

$$\begin{aligned}\hat{y}_i(t-1) &= w_0(t-1) + w_1(t-1)x_i, \quad \forall i, \\ e_i(t-1) &= y_i - \hat{y}_i(t-1), \quad \forall i.\end{aligned}$$

- 3 Atualize os parâmetros:

$$\begin{aligned}w_0(t) &= w_0(t-1) + \alpha \frac{1}{N} \sum_{i=1}^N e_i(t-1) \\ w_1(t) &= w_1(t-1) + \alpha \frac{1}{N} \sum_{i=1}^N e_i(t-1)x_i\end{aligned}$$

Regressão linear simples - Otimização via GD



Regressão linear simples

Gradiente Descendente Estocástico (SGD, stochastic gradient descent)

- 1 Escolha um valor α positivo e pequeno.
- 2 Inicialize os parâmetros do modelo na iteração $t = 0$.
- 3 Repita por diversos ciclos (épocas):
 - 1 Permute aleatoriamente a ordem dos dados.
 - 2 Para cada padrão de entrada, $i = 1, \dots, N$, repita:
 - 1 Faça $t \leftarrow t + 1$.
 - 2 Calcule os erros do modelo:

$$\hat{y}_i(t-1) = w_0(t-1) + w_1(t-1)x_i,$$

$$e_i(t-1) = y_i - \hat{y}_i(t-1).$$

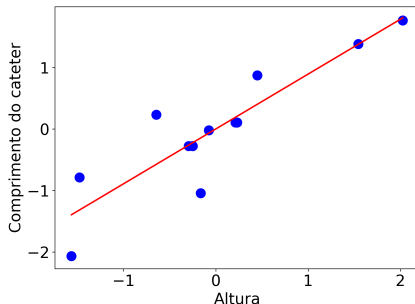
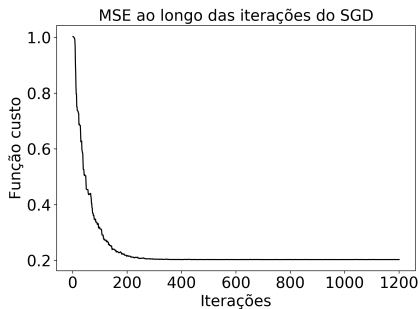
- 3 Atualize os parâmetros:

$$w_0(t) = w_0(t-1) + \alpha e_i(t-1)$$

$$w_1(t) = w_1(t-1) + \alpha e_i(t-1)x_i$$

- Também chamado de **algoritmo LMS** (*Least Mean Squares*).

Regressão linear simples - Otimização via SGD



Regressão linear múltipla

- Podemos reconsiderar o problema inserindo o peso do paciente:

Altura (m)	Peso (Kg)	Comprimento (cm)
1.087	18.141	37
1.613	42.404	50
0.953	16.100	34
1.003	13.605	36
1.156	23.583	43
0.978	7.710	28
1.092	17.460	37
0.572	3.855	20
0.940	14.966	34
0.597	4.308	30
0.838	9.524	38
1.473	35.828	47

Regressão linear múltipla

- Podemos reconsiderar o problema inserindo o peso do paciente:

Altura (m)	Peso (Kg)	Comprimento (cm)
1.087	18.141	37
1.613	42.404	50
0.953	16.100	34
1.003	13.605	36
1.156	23.583	43
0.978	7.710	28
1.092	17.460	37
0.572	3.855	20
0.940	14.966	34
0.597	4.308	30
0.838	9.524	38
1.473	35.828	47

- Novo modelo linear **múltiplo**:

$$\hat{y}_i = w_0 + w_1 x_{i1} + w_2 x_{i2}$$

- x_{i1} é a i -ésima **Altura** e x_{i2} é o i -ésimo **Peso**.

Regressão linear múltipla

- Caso façamos $x_{i0} = 1$, temos:

$$\hat{y}_i = w_0 x_{i0} + w_1 x_{i1} + w_2 x_{i2}$$

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i.$$

- Note que:

$$\mathbf{w} = [w_0, w_1, w_2]^\top,$$

$$\mathbf{x}_i = [1, x_{i1}, x_{i2}]^\top.$$

Regressão linear múltipla

Gradiente Descendente

- Regra de atualização:

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \alpha \frac{1}{N} \sum_{i=1}^N e_i(t-1) \mathbf{x}_i$$

Gradiente Descendente Estocástico

- Regra de atualização:

$$\mathbf{w}(t) = \mathbf{w}(t-1) + \alpha e_i(t-1) \mathbf{x}_i$$

Diferença probabilística entre o GD e o SGD

- O erro quadrático \mathcal{J} é uma **variável aleatória**:

$$\min_{\mathbf{w}} \frac{1}{2} \mathbb{E}\{\mathcal{J}(\mathbf{w})\}, \quad \mathcal{J}(\mathbf{w}) = e^2 = (y - \hat{y})^2,$$

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} \mathbb{E}\{\mathcal{J}\}$$

Diferença probabilística entre o GD e o SGD

- O erro quadrático \mathcal{J} é uma **variável aleatória**:

$$\min_{\mathbf{w}} \frac{1}{2} \mathbb{E}\{\mathcal{J}(\mathbf{w})\}, \quad \mathcal{J}(\mathbf{w}) = e^2 = (y - \hat{y})^2,$$

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} \mathbb{E}\{\mathcal{J}\}$$

- A **média amostral** resulta no algoritmo GD:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right)$$

Diferença probabilística entre o GD e o SGD

- O erro quadrático \mathcal{J} é uma **variável aleatória**:

$$\min_{\mathbf{w}} \frac{1}{2} \mathbb{E}\{\mathcal{J}(\mathbf{w})\}, \quad \mathcal{J}(\mathbf{w}) = e^2 = (y - \hat{y})^2,$$

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} \mathbb{E}\{\mathcal{J}\}$$

- A **média amostral** resulta no algoritmo GD:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right)$$

- Uma **aproximação estocástica** resulta no algoritmo SGD:

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{2} \frac{\partial}{\partial \mathbf{w}} (y_i - \hat{y}_i)^2$$

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

Regressão linear analítica

- Reunimos todos os padrões de entrada \mathbf{x}_i em uma matriz \mathbf{X} :

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_N]^\top \in \mathbb{R}^{N \times (D+1)}.$$

- N é o número de observações/amostras/vetores/padrões.
- D é a dimensão da entrada (excluindo o termo $x_{i0} = 1$).

- Agrupamos as saídas disponíveis em um vetor \mathbf{y} :

$$\mathbf{y} = [y_1, y_2, \cdots, y_N]^\top \in \mathbb{R}^N.$$

- Agrupamos as saídas do modelo em um vetor $\hat{\mathbf{y}}$:

$$\hat{\mathbf{y}} = [\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_N]^\top \in \mathbb{R}^N.$$

- Agrupamos os parâmetros em um vetor \mathbf{w} :

$$\mathbf{w} = [w_0, w_1, \cdots, w_D]^\top \in \mathbb{R}^{D+1}.$$

Regressão linear analítica

- Dados do problema do cateterismo cardíaco em formato matricial:

$$\mathbf{X} = \begin{bmatrix} 1 & 1.087 & 18.141 \\ 1 & 1.613 & 42.404 \\ 1 & 0.953 & 16.100 \\ 1 & 1.003 & 13.605 \\ 1 & 1.156 & 23.583 \\ 1 & 0.978 & 7.710 \\ 1 & 1.092 & 17.460 \\ 1 & 0.572 & 3.855 \\ 1 & 0.940 & 14.966 \\ 1 & 0.597 & 4.308 \\ 1 & 0.838 & 9.524 \\ 1 & 1.473 & 35.828 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 37 \\ 50 \\ 34 \\ 36 \\ 43 \\ 28 \\ 37 \\ 20 \\ 34 \\ 30 \\ 38 \\ 47 \end{bmatrix}$$

Regressão linear analítica

- Reformulamos nosso modelo linear na forma matricial:

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i,$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}.$$

Regressão linear analítica

- Reformulamos nosso modelo linear na forma matricial:

$$\hat{y}_i = \mathbf{w}^\top \mathbf{x}_i,$$

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}.$$

- Reformulamos também a função custo:

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

Regressão linear analítica

- O mínimo de $\mathcal{J}(\mathbf{w})$ ocorrerá em $\frac{\partial \mathcal{J}(\mathbf{w})}{\partial \mathbf{w}} = 0$:

$$\begin{aligned}\frac{\partial \mathcal{J}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{1}{2} 2(-\mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w}) \\ -\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X}\mathbf{w} &= 0 \\ \mathbf{X}^\top \mathbf{X}\mathbf{w} &= \mathbf{X}^\top \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

Regressão linear analítica

- O mínimo de $\mathcal{J}(\mathbf{w})$ ocorrerá em $\frac{\partial \mathcal{J}(\mathbf{w})}{\partial \mathbf{w}} = 0$:

$$\begin{aligned}\frac{\partial \mathcal{J}(\mathbf{w})}{\partial \mathbf{w}} &= \frac{1}{2} 2(-\mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\mathbf{w}) \\ -\mathbf{X}^\top \mathbf{y} + \mathbf{X}^\top \mathbf{X}\mathbf{w} &= 0 \\ \mathbf{X}^\top \mathbf{X}\mathbf{w} &= \mathbf{X}^\top \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

- $\mathbf{X}^+ = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$, em que $\mathbf{X}^+ \mathbf{X} = \mathbf{I}$, é chamada de **inversa de Moore-Penrose** ou **pseudo-inversa**.

Regressão linear analítica

Método dos mínimos quadrados ordinários (OLS, *ordinary least squares*)

- O vetor de parâmetros \mathbf{w} que minimiza $\mathcal{J}(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$ é dado por

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y},$$

em quem \mathbf{X} é a matriz de vetores de entrada (um por linha) e \mathbf{y} é o vetor de saídas desejadas.

Regressão linear analítica

- OLS equivale ao método de Newton aplicado na função de custo quadrática \mathcal{J} :

$$\begin{aligned} \mathbf{w} &= \mathbf{w}_0 - \left(\frac{\partial^2 \mathcal{J}(\mathbf{w}_0)}{\partial \mathbf{w}_0^2} \right)^{-1} \frac{\partial \mathcal{J}(\mathbf{w}_0)}{\partial \mathbf{w}_0}, \\ \frac{\partial \mathcal{J}(\mathbf{w}_0)}{\partial \mathbf{w}_0} &= -\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}_0), \\ \frac{\partial^2 \mathcal{J}(\mathbf{w}_0)}{\partial \mathbf{w}_0^2} &= \mathbf{X}^\top \mathbf{X}, \\ \mathbf{w} &= \mathbf{w}_0 - (\mathbf{X}^\top \mathbf{X})^{-1} (-\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}_0)) \end{aligned}$$

- Podemos escolher $\mathbf{w}_0 = \mathbf{0}$ para obter o mínimo global.

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

Regressão linear analítica

- De onde vem a função custo $\mathcal{J}(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$?

Regressão linear analítica

- De onde vem a função custo $\mathcal{J}(\mathbf{w}) = \frac{1}{2}(\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$?
- Considerando um ruído independente $\epsilon \sim \mathcal{N}(0, \sigma^2)$:

$$y_i = \hat{y}_i + \epsilon = \mathbf{w}^\top \mathbf{x}_i + \epsilon,$$

$$p(y_i | \mathbf{x}_i, \mathbf{w}) = \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2),$$

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \prod_{i=1}^N \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2)$$

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \sum_{i=1}^N \log \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left(-\frac{(y_i - \mathbf{w}^\top \mathbf{x}_i)^2}{2\sigma^2} \right)$$

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \underbrace{-\frac{N}{2} \log(2\pi\sigma^2)}_{\text{const. em relação a } \mathbf{w}} - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

Regressão linear analítica

- Queremos maximizar $\mathcal{L}(\mathbf{w}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w})$, o que equivale a minimizar $\mathcal{J}(\mathbf{w}) = -\mathcal{L}(\mathbf{w})$.
- Ignorando os termos constantes:

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$$

Regressão linear analítica

- Queremos maximizar $\mathcal{L}(\mathbf{w}) = \log p(\mathbf{y}|\mathbf{X}, \mathbf{w})$, o que equivale a minimizar $\mathcal{J}(\mathbf{w}) = -\mathcal{L}(\mathbf{w})$.
- Ignorando os termos constantes:

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \mathbf{w}^\top \mathbf{x}_i)^2$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

$$\mathcal{J}(\mathbf{w}) = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$$

Solução de máxima verossimilhança

A solução obtida via OLS (e aproximada via GD e SGD), chamada de solução de **máxima verossimilhança (maximum likelihood)**, é ótima quando o ruído é Gaussiano:

$$\mathbf{w}_{\text{OLS}} = \mathbf{w}_{\text{ML}} = \arg \max \log p(\mathbf{y}|\mathbf{X}, \mathbf{w})$$

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

Propriedades da Distribuição Gaussiana

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}.$$

Marginalização

A observação de uma coleção maior de variáveis não afeta a distribuição de subconjuntos menores, ou seja:

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}) \text{ e } \mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$$

Condicionamento

Condicionar Gaussianas resulta em uma Gaussiana:

$$p(\mathbf{x}_1 | \mathbf{x}_2 = \mathbf{z}) = \mathcal{N}(\mathbf{x}_1 | \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{z} - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21})$$

Linearidade

Uma combinação linear de Gaussianas resulta em uma Gaussiana:

$$p(a\mathbf{x}_1 + b\mathbf{x}_2) = \mathcal{N}(a\boldsymbol{\mu}_1 + b\boldsymbol{\mu}_2, a^2\boldsymbol{\Sigma}_{11} + b^2\boldsymbol{\Sigma}_{22})$$

Regressão Linear Bayesiana

- Considere um modelo linear com ruído Gaussiano:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \quad \epsilon_i \sim \mathcal{N}(\epsilon_i|0, \sigma^2).$$

Regressão Linear Bayesiana

- Considere um modelo linear com ruído Gaussiano:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \quad \epsilon_i \sim \mathcal{N}(\epsilon_i|0, \sigma^2).$$

- A verossimilhança do modelo é dada por:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}).$$

Regressão Linear Bayesiana

- Considere um modelo linear com ruído Gaussiano:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \quad \epsilon_i \sim \mathcal{N}(\epsilon_i|0, \sigma^2).$$

- A verossimilhança do modelo é dada por:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}).$$

- Por ser desconhecido, o vetor de parâmetros \mathbf{w} recebe uma distribuição a priori, por exemplo, Gaussiana:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$

Regressão Linear Bayesiana

- Considere um modelo linear com ruído Gaussiano:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}, \quad \epsilon_i \sim \mathcal{N}(\epsilon_i|0, \sigma^2).$$

- A verossimilhança do modelo é dada por:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I}).$$

- Por ser desconhecido, o vetor de parâmetros \mathbf{w} recebe uma distribuição a priori, por exemplo, Gaussiana:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0).$$

- Após a observação dos dados $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, como calcular a distribuição a posteriori $p(\mathbf{w}|\mathcal{D})$?

Regressão Linear Bayesiana

- Voltamos às propriedades da Gaussiana (com nova notação):

$$\mathbf{z} = \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}(\mathbf{m}_z, \mathbf{S}_z), \quad \mathbf{m}_z = \begin{bmatrix} \boldsymbol{\mu}_w \\ \boldsymbol{\mu}_y \end{bmatrix}, \quad \mathbf{S}_z = \begin{bmatrix} \boldsymbol{\Sigma}_w & \boldsymbol{\Sigma}_{wy} \\ \boldsymbol{\Sigma}_{yw} & \boldsymbol{\Sigma}_y \end{bmatrix}.$$

Marginalização

A observação de uma coleção maior de variáveis não afeta a distribuição de subconjuntos menores, ou seja:

$$\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w) \text{ e } \mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$$

Condicionamento

Condicionar Gaussianas resulta em uma Gaussiana:

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}_w + \boldsymbol{\Sigma}_{wy}\boldsymbol{\Sigma}_y^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \boldsymbol{\Sigma}_w - \boldsymbol{\Sigma}_{wy}\boldsymbol{\Sigma}_y^{-1}\boldsymbol{\Sigma}_{yw})$$

Regressão Linear Bayesiana

- Usamos a priori de \mathbf{w} , $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$:

$$\mathbf{z} = \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}(\mathbf{m}_z, \mathbf{S}_z), \quad \mathbf{m}_z = \begin{bmatrix} \mathbf{m}_0 \\ \mu_y \end{bmatrix}, \quad \mathbf{S}_z = \begin{bmatrix} \mathbf{S}_0 & \Sigma_{wy} \\ \Sigma_{yw} & \Sigma_y \end{bmatrix}.$$

Marginalização

A observação de uma coleção maior de variáveis não afeta a distribuição de subconjuntos menores, ou seja:

$$\mathbf{w} \sim \mathcal{N}(\mathbf{m}_0, \mathbf{S}_0) \text{ e } \mathbf{y} \sim \mathcal{N}(\mu_y, \Sigma_y)$$

Condicionamento

Condicionar Gaussianas resulta em uma Gaussiana:

$$p(\mathbf{w} | \mathbf{y}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0 + \Sigma_{wy} \Sigma_y^{-1}(\mathbf{y} - \mu_y), \mathbf{S}_0 - \Sigma_{wy} \Sigma_y^{-1} \Sigma_{yw})$$

Regressão Linear Bayesiana

- Note que pela propriedade do condicionamento, podemos obter a posteriori de \mathbf{w} :

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0 + \Sigma_{wy} \Sigma_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y), \mathbf{S}_0 - \Sigma_{wy} \Sigma_y^{-1} \Sigma_{yw}).$$

Regressão Linear Bayesiana

- Note que pela propriedade do condicionamento, podemos obter a posteriori de \mathbf{w} :

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0 + \Sigma_{wy} \Sigma_y^{-1} (\mathbf{y} - \boldsymbol{\mu}_y), \mathbf{S}_0 - \Sigma_{wy} \Sigma_y^{-1} \Sigma_{yw}).$$

- Para calcular a expressão acima, precisamos de alguns valores ainda desconhecidos:
 - $\boldsymbol{\mu}_y$ e Σ_y , provenientes da verossimilhança marginal $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \Sigma_y)$.
 - Σ_{wy} , a matriz de covariância cruzada entre \mathbf{w} e \mathbf{y} .

Regressão Linear Bayesiana

- Começamos calculando a verossimilhança marginal de \mathbf{y} :

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)p(\mathbf{w})d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I}). \end{aligned}$$

Regressão Linear Bayesiana

- Começamos calculando a verossimilhança marginal de \mathbf{y} :

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2)p(\mathbf{w})d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I}). \end{aligned}$$

- Portanto, se $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$:

$$\begin{aligned} \boldsymbol{\mu}_y &= \mathbf{X}\mathbf{m}_0, \\ \boldsymbol{\Sigma}_y &= \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I}. \end{aligned}$$

Regressão Linear Bayesiana

- Começamos calculando a verossimilhança marginal de \mathbf{y} :

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \sigma^2) p(\mathbf{w}) d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{m}_0, \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}). \end{aligned}$$

- Portanto, se $\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_y, \boldsymbol{\Sigma}_y)$:

$$\begin{aligned} \boldsymbol{\mu}_y &= \mathbf{X}\mathbf{m}_0, \\ \boldsymbol{\Sigma}_y &= \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2 \mathbf{I}. \end{aligned}$$

- Qual a covariância entre \mathbf{w} e \mathbf{y} , isto é, $\boldsymbol{\Sigma}_{wy}$?

$$\begin{aligned} \boldsymbol{\Sigma}_{wy} &= \text{cov}[\mathbf{w}, \mathbf{y}] = \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon}] \\ &= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w}] + \text{cov}[\mathbf{w}, \boldsymbol{\epsilon}] \\ &= \text{cov}[\mathbf{w}, \mathbf{X}\mathbf{w}] = \text{cov}[\mathbf{w}, \mathbf{w}] \mathbf{X}^\top = \mathbf{S}_0 \mathbf{X}^\top. \end{aligned}$$

Regressão Linear Bayesiana

- Agora podemos construir a distribuição conjunta $p(\mathbf{w}, \mathbf{y})$:

$$\mathbf{z} = \begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}(\mathbf{m}_z, \mathbf{S}_z), \quad \mathbf{m}_z = \begin{bmatrix} \mathbf{m}_0 \\ \boldsymbol{\mu}_y \end{bmatrix}, \quad \mathbf{S}_z = \begin{bmatrix} \mathbf{S}_0 & \boldsymbol{\Sigma}_{wy} \\ \boldsymbol{\Sigma}_{yw} & \boldsymbol{\Sigma}_y \end{bmatrix},$$

$$\begin{bmatrix} \mathbf{w} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_0 \\ \mathbf{X}\mathbf{m}_0 \end{bmatrix}, \begin{bmatrix} \mathbf{S}_0 & \mathbf{S}_0\mathbf{X}^\top \\ \mathbf{X}\mathbf{S}_0 & \mathbf{X}\mathbf{S}_0\mathbf{X}^\top + \sigma^2\mathbf{I} \end{bmatrix} \right).$$

Regressão Linear Bayesiana

- Finalmente, usamos a propriedade de condicionamento de Gaussianas para obter a posteriori de \mathbf{w} :

$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0 + \Sigma_{wy}\Sigma_y^{-1}(\mathbf{y} - \boldsymbol{\mu}_y), \mathbf{S}_0 - \Sigma_{wy}\Sigma_y^{-1}\Sigma_{yw}),$$

$$p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

$$\boldsymbol{\mu} = \mathbf{m}_0 + \mathbf{S}_0\mathbf{X}^\top(\mathbf{XS}_0\mathbf{X}^\top + \sigma^2\mathbf{I})^{-1}(\mathbf{y} - \mathbf{Xm}_0)$$

$$= \mathbf{m}_0 + (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top(\mathbf{y} - \mathbf{Xm}_0),$$

$$\boldsymbol{\Sigma} = \mathbf{S}_0 - \mathbf{S}_0\mathbf{X}^\top(\mathbf{XS}_0\mathbf{X}^\top + \sigma^2\mathbf{I})^{-1}\mathbf{XS}_0$$

$$= \mathbf{S}_0 - (\mathbf{S}_0\mathbf{X}^\top\mathbf{X} + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\mathbf{X}^\top\mathbf{XS}_0.$$

em que foi usada a identidade matricial abaixo (derivada a partir da identidade de Woodbury):

$$\mathbf{U}(\mathbf{VU} + \mathbf{I})^{-1} = (\mathbf{UV} + \mathbf{I})^{-1}\mathbf{U}$$

- Note que o modelo não é definido por um único valor para \mathbf{w} , mas uma distribuição a posteriori $p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2)$.

Regressão Linear Bayesiana

- Com a distribuição a posteriori $p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2)$, podemos fazer previsões para novas entradas \mathbf{X}_* :

$$\mathbf{y}_* = \mathbf{X}_* \mathbf{w} + \epsilon,$$

$$\begin{aligned} p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \int p(\mathbf{y}_*|\mathbf{w}, \sigma^2) p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) d\mathbf{w} \\ &= \int \mathcal{N}(\mathbf{y}_*|\mathbf{X}_* \mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}_*|\mathbf{X}_* \boldsymbol{\mu}, \mathbf{X}_* \boldsymbol{\Sigma} \mathbf{X}_*^\top + \sigma^2 \mathbf{I}). \end{aligned}$$

- Note que a predição não é um único valor, mas uma distribuição de probabilidade bem definida.

Regressão Linear Bayesiana

- Caso escolhamos uma priori $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_w^2 \mathbf{I})$, temos:

$$\begin{aligned}p(\mathbf{w}|\mathcal{D}, \sigma_w^2, \sigma^2) &= \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ \boldsymbol{\mu} &= \sigma_w^2 (\sigma_w^2 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}, \\ \boldsymbol{\Sigma} &= \sigma_w^2 \mathbf{I} - \sigma_w^2 (\sigma_w^2 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{X} \sigma_w^2.\end{aligned}$$

- Podemos reescrever a média a posteriori $\boldsymbol{\mu}$ para obter a solução de mínimos quadrados regularizada (ridge regression):

$$\begin{aligned}\boldsymbol{\mu} &= \left(\mathbf{X}^\top \mathbf{X} + \frac{\sigma^2}{\sigma_w^2} \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y} \\ &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}, \quad \lambda = \frac{\sigma^2}{\sigma_w^2}.\end{aligned}$$

Regressão Linear Bayesiana

Resumo do algoritmo

- Passo de estimação

- 1 Defina a partir de conhecimentos/experimentos anteriores:

- os momentos da priori $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$;

- a variância do ruído $p(\epsilon) = \mathcal{N}(\epsilon|0, \sigma^2)$.

- 2 A partir dos dados $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, calcule a posteriori de \mathbf{w} :

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

$$\boldsymbol{\mu} = \mathbf{m}_0 + (\mathbf{S}_0 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \mathbf{m}_0),$$

$$\boldsymbol{\Sigma} = \mathbf{S}_0 - (\mathbf{S}_0 \mathbf{X}^\top \mathbf{X} + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \mathbf{X}^\top \mathbf{X} \mathbf{S}_0.$$

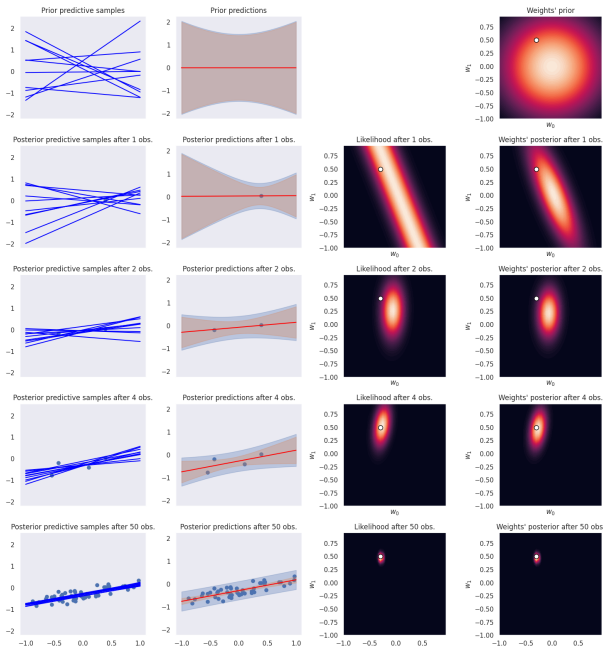
- 3 Retorne a posteriori $p(\mathbf{w}|\mathcal{D})$ dos parâmetros.

- Passo de predição

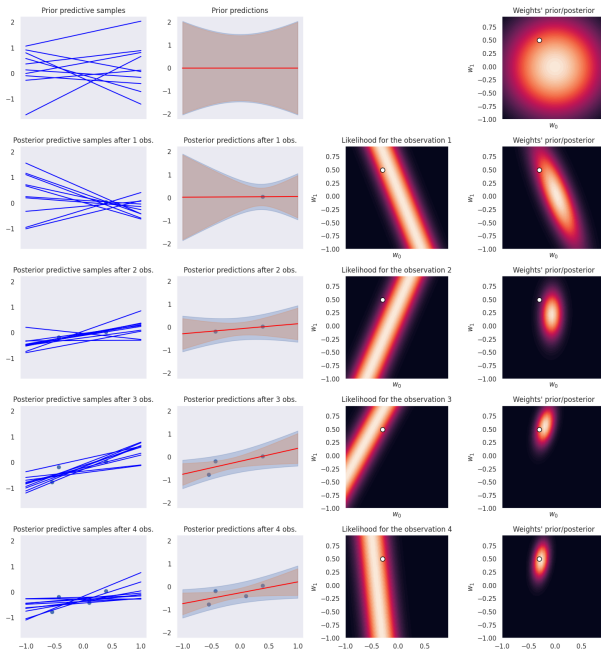
- 1 Dados padrões $\mathbf{X}_* \in \mathbb{R}^{(N_* \times D)}$, retorne a distribuição preditiva

$$p(\mathbf{y}_*|\mathbf{X}_*) = \mathcal{N}(\mathbf{y}_*|\mathbf{X}_* \boldsymbol{\mu}, \mathbf{X}_* \boldsymbol{\Sigma} \mathbf{X}_*^\top + \sigma^2 \mathbf{I}).$$

Regressão Linear Bayesiana em batch



Regressão Linear Bayesiana sequencial



Regressão Linear Bayesiana com Funções de Base

- Podemos usar uma transformação arbitrária $\phi(\cdot)$ nos atributos:

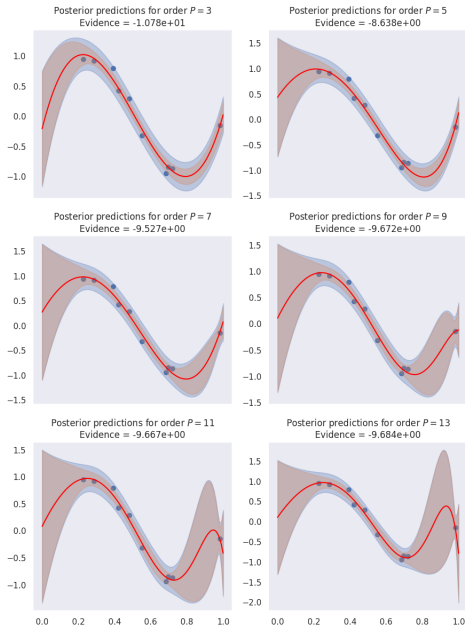
$$\Phi = \phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \phi(\mathbf{x}_2), \dots, \phi(\mathbf{x}_N)]^\top, \quad \phi: \mathbb{R}^D \rightarrow \mathbb{R}^Q,$$
$$\mathbf{y} = \Phi \mathbf{w} + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon|0, \sigma^2).$$

- As expressões do modelo linear Bayesiano continuam as mesmas:

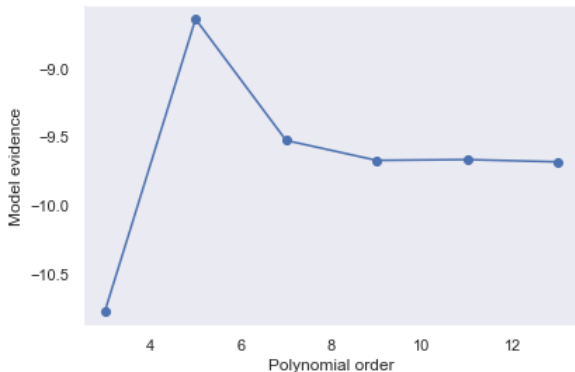
$$p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\Phi \mathbf{w}, \sigma^2 \mathbf{I}),$$
$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0),$$
$$p(\mathbf{w}|\mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma}), \quad \mathcal{D} = (\Phi, \mathbf{y}),$$
$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{y}_*|\Phi_* \boldsymbol{\mu}, \Phi_* \boldsymbol{\Sigma} \Phi_*^\top + \sigma^2 \mathbf{I}),$$
$$\boldsymbol{\mu} = \mathbf{m}_0 + (\mathbf{S}_0 \Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \Phi^\top (\mathbf{y} - \Phi \mathbf{m}_0),$$
$$\boldsymbol{\Sigma} = \mathbf{S}_0 - (\mathbf{S}_0 \Phi^\top \Phi + \sigma^2 \mathbf{I})^{-1} \mathbf{S}_0 \Phi^\top \Phi \mathbf{S}_0.$$

- Regressão polinomial Bayesiana de ordem P pode ser obtida fazendo $\phi(x_i) = [1, x_i, x_i^2, \dots, x_i^P]^\top$.

Regressão Polinomial Bayesiana variando a ordem



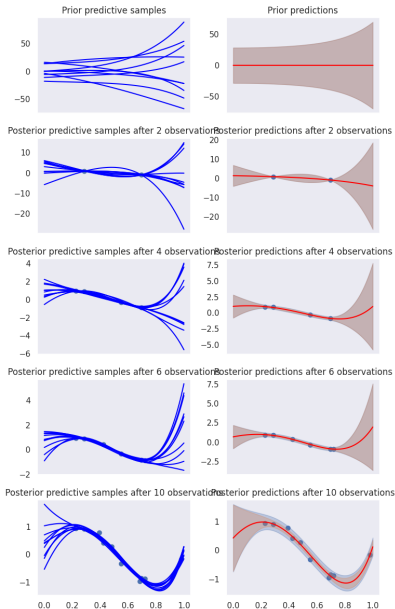
Regressão Polinomial Bayesiana variando a ordem



Lembrando que a evidência do modelo é dada por:

$$\begin{aligned}\log p(\mathbf{y}|\Phi, \mathbf{m}_0, \mathbf{S}_0, \sigma^2) &= \log \int p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2) p(\mathbf{w}) d\mathbf{w} \\ &= \log \int \mathcal{N}(\mathbf{y}|\Phi \mathbf{w}, \sigma^2 \mathbf{I}) \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0) d\mathbf{w} \\ &= \log \mathcal{N}(\mathbf{y}|\Phi \mathbf{m}_0, \Phi \mathbf{S}_0 \Phi^\top + \sigma^2 \mathbf{I}).\end{aligned}$$

Regressão Polinomial Bayesiana em batch



Regressão Linear Bayesiana com Funções de Base

- Podemos escolher $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{S}_0)$:

$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{y}_*|\Phi_*\boldsymbol{\mu}, \Phi_*\Sigma\Phi_*^\top + \sigma^2\mathbf{I}),$$

$$\boldsymbol{\mu} = \underbrace{(\mathbf{S}_0\Phi^\top\Phi + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\Phi^\top}_{\mathbf{A}}\mathbf{y},$$

$$\Sigma = \mathbf{S}_0 - \underbrace{(\mathbf{S}_0\Phi^\top\Phi + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\Phi^\top}_{\mathbf{A}}\Phi\mathbf{S}_0.$$

- Usamos a seguinte identidade matricial (derivada da identidade de Woodbury):

$$\mathbf{A} = (\mathbf{S}_0\Phi^\top\Phi + \sigma^2\mathbf{I})^{-1}\mathbf{S}_0\Phi^\top = \mathbf{S}_0\Phi^\top(\Phi\mathbf{S}_0\Phi^\top + \sigma^2\mathbf{I})^{-1}$$

- Reescrevemos a distribuição preditiva:

$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{y}_*|\Phi_*\mathbf{S}_0\Phi^\top(\Phi\mathbf{S}_0\Phi^\top + \sigma^2\mathbf{I})^{-1}\mathbf{y},$$

$$\Phi_*(\mathbf{S}_0 - \mathbf{S}_0\Phi^\top(\Phi\mathbf{S}_0\Phi^\top + \sigma^2\mathbf{I})^{-1}\Phi\mathbf{S}_0)\Phi_*^\top + \sigma^2\mathbf{I}),$$

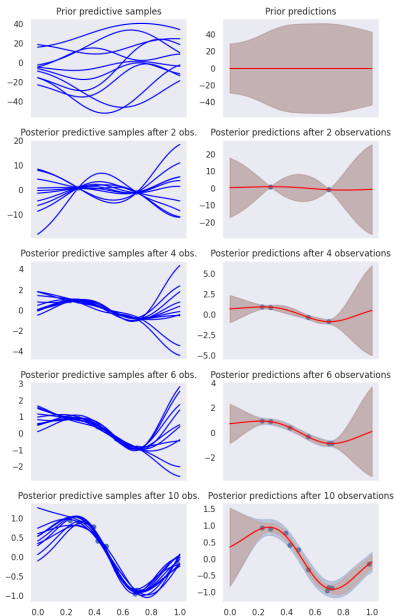
Regressão RBF Bayesiana

- Podemos usar uma função de base radial (radial basis function) para obter uma regressão RBF Bayesiana.
- Considerando M funções RBF com centros $c_m|_{m=1}^M$, largura de banda $\lambda > 0$ e entradas unidimensionais:

$$\begin{aligned}y_i &= \mathbf{w}^\top \boldsymbol{\phi}_i + \epsilon, \\ \boldsymbol{\phi}_i &= [1, \phi_1(x_i), \dots, \phi_M(x_i)]^\top, \\ \phi_m(x_i) &= \exp\left(-\frac{(x_i - c_m)^2}{2\lambda}\right).\end{aligned}$$

- Por exemplo, podemos inicializar os centros via algoritmo de agrupamento e fazer $\lambda = \mathbb{V}(x)$.

Regressão RBF Bayesiana



Do espaço de atributos para funções de kernel

- Fazendo $\Psi = \Phi S_0^{1/2}$ e $\Psi_* = \Phi_* S_0^{1/2}$ a preditiva se torna:

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathcal{D}, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{y}_* | \Psi_* \Psi^\top (\Psi \Psi^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \\ (\Psi_* \Psi_*^\top - \Psi_* \Psi^\top (\Psi \Psi^\top + \sigma^2 \mathbf{I})^{-1} \Psi \Psi_*^\top) + \sigma^2 \mathbf{I}),$$

- Agora podemos aplicar o *kernel trick* (truque do kernel):

$$\begin{aligned} \Psi \Psi^\top &= k(\mathbf{X}, \mathbf{X}) = \mathbf{K}, \\ \Psi_* \Psi^\top &= k(\mathbf{X}_*, \mathbf{X}) = \mathbf{K}_{*f}, \\ \Psi \Psi_*^\top &= k(\mathbf{X}, \mathbf{X}_*) = \mathbf{K}_{f*}, \\ \Psi_* \Psi_*^\top &= k(\mathbf{X}_*, \mathbf{X}_*) = \mathbf{K}_{**}. \end{aligned}$$

- Finalmente, obtemos a versão abaixo da distribuição preditiva:

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathcal{D}, \mathbf{S}_0, \sigma^2) = \mathcal{N}(\mathbf{y}_* | \mathbf{K}_{*f} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \\ \mathbf{K}_{**} - \mathbf{K}_{*f} (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{f*} + \sigma^2 \mathbf{I}),$$

correspondente a um modelo de processo Gaussiano.

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

Comparação Bayesiana de modelos

- A abordagem Bayesiana evita otimizar parâmetros, preferindo marginalizá-los.
- Não há a necessidade de um conjunto de validação para comparar modelos.
- Dado um modelo \mathcal{M}_j e um conjunto de dados $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, desejamos computar a distribuição a posteriori do modelo:

$$p(\mathcal{M}_j|\mathcal{D}) = \frac{p(\mathcal{M}_j)p(\mathcal{D}|\mathcal{M}_j)}{p(\mathcal{D})} \propto p(\mathcal{M}_j)p(\mathcal{D}|\mathcal{M}_j).$$

Comparação Bayesiana de modelos

- A abordagem Bayesiana evita otimizar parâmetros, preferindo marginalizá-los.
- Não há a necessidade de um conjunto de validação para comparar modelos.
- Dado um modelo \mathcal{M}_j e um conjunto de dados $\mathcal{D} = (\mathbf{X}, \mathbf{y})$, desejamos computar a distribuição a posteriori do modelo:

$$p(\mathcal{M}_j|\mathcal{D}) = \frac{p(\mathcal{M}_j)p(\mathcal{D}|\mathcal{M}_j)}{p(\mathcal{D})} \propto p(\mathcal{M}_j)p(\mathcal{D}|\mathcal{M}_j).$$

- Considere L modelos com a priori $p(\mathcal{M}_j)$ iguais.
- Evidência do modelo \mathcal{M}_j : $\propto p(\mathcal{D}|\mathcal{M}_j)$.
- **Fator de Bayes** entre modelos com priori iguais:

$$B_{j,k} = \frac{p(\mathcal{D}|\mathcal{M}_j)}{p(\mathcal{D}|\mathcal{M}_k)}.$$

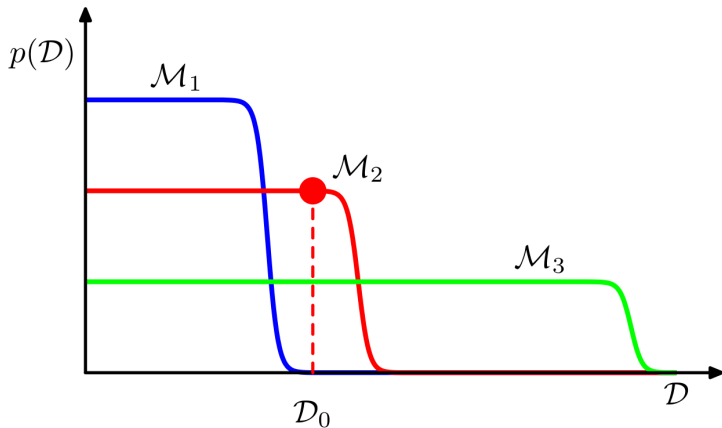
Evidência e Navalha de Occam Bayesiana

- A evidência é obtida pela marginalização dos parâmetros:

$$p(\mathcal{D}|\mathcal{M}_j) = \int p(\mathcal{D}|\mathbf{w}, \mathcal{M}_j)p(\mathbf{w}|\mathcal{M}_j)d\mathbf{w}$$
$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w} \quad (\text{omitindo } \mathcal{M}_j).$$

- Soluções ML ou MAP usam somente o ajuste aos dados $p(\mathcal{D}|\mathbf{w})$ na comparação de modelos, podendo resultar em overfitting.
- Ao marginalizar os parâmetros, penalizamos modelos muito complexos.
- **Navalha de Occam Bayesiana:** Equilíbrio entre ajuste aos dados e complexidade na comparação Bayesiana de modelos.

Conservação de massa de probabilidade



Mistura e seleção de modelos

- **Mistura de modelos:** Predições são feitas combinando todos os modelos (note que os parâmetros já foram marginalizados):

$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}) = \sum_{j=1}^L p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{M}_j, \mathcal{D})p(\mathcal{M}_j|\mathcal{D}).$$

Mistura e seleção de modelos

- **Mistura de modelos:** Predições são feitas combinando todos os modelos (note que os parâmetros já foram marginalizados):

$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}) = \sum_{j=1}^L p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{M}_j, \mathcal{D})p(\mathcal{M}_j|\mathcal{D}).$$

- **Seleção de modelos:** Somente o modelo com maior posteriori $p(\mathcal{M}_{\text{best}}|\mathcal{D})$ é usado na predição.
 - Para modelos com priori iguais, equivale a usar aquele com maior verossimilhança marginal $p(\mathcal{D}|\mathcal{M}_{\text{best}})$.

Evidência aproximada

- Em um tratamento Bayesiano puro, colocaríamos priori nos hiperparâmetros φ (por exemplo σ_w^2 e σ^2), que seriam marginalizados junto com os parâmetros para fazer predições:

$$p(\mathbf{y}_* | \mathbf{X}_*, \mathcal{D}) = \int p(\mathbf{y}_* | \mathbf{X}_*, \mathcal{D}, \mathbf{w}, \varphi) p(\mathbf{w} | \mathcal{D}, \varphi) p(\varphi | \mathcal{D}) d\mathbf{w} d\varphi.$$

Evidência aproximada

- A integral anterior usualmente não é analítica, podendo ser aproximada por $p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}) \approx p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \hat{\varphi})$:

$$p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \hat{\varphi}) = \int p(\mathbf{y}_*|\mathbf{X}_*, \mathcal{D}, \mathbf{w}, \hat{\varphi})p(\mathbf{w}|\mathcal{D}, \hat{\varphi})d\mathbf{w},$$

$$\hat{\varphi} = \arg \max p(\mathcal{D}|\varphi) = \arg \max \int p(\mathcal{D}|\mathbf{w}, \varphi)p(\mathbf{w}|\varphi)d\mathbf{w}.$$

- O procedimento acima também é chamado de *empirical Bayes*, *type 2 maximum likelihood* ou *generalized maximum likelihood*.

Computando a evidência

- Considere o modelo de regressão abaixo:

$$p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}), \quad \Phi = \phi(\mathbf{X}),$$
$$p(\mathbf{w}|\sigma_w^2) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_w^2\mathbf{I}).$$

Computando a evidência

- Considere o modelo de regressão abaixo:

$$p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}), \quad \Phi = \phi(\mathbf{X}),$$
$$p(\mathbf{w}|\sigma_w^2) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_w^2\mathbf{I}).$$

- A evidência é dada por:

$$p(\mathbf{y}|\Phi) = \int p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2)p(\mathbf{w}|\sigma_w^2)p(\sigma^2, \sigma_w^2)d\mathbf{w}d\sigma_w^2d\sigma^2$$
$$\approx p(\mathbf{y}|\Phi, \sigma^2, \sigma_w^2) = \int p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2)p(\mathbf{w}|\sigma_w^2)d\mathbf{w},$$
$$p(\mathbf{y}|\Phi, \sigma^2, \sigma_w^2) = \int \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I})\mathcal{N}(\mathbf{w}|\mathbf{0}, \sigma_w^2\mathbf{I})d\mathbf{w}$$
$$= \mathcal{N}(\mathbf{y}|\mathbf{0}, \sigma_w^2\Phi\Phi^\top + \sigma^2\mathbf{I}).$$

- A otimização pode ser feita iterativamente a partir dos dados de treinamento e de valores iniciais (Bishop, pg. 168).

Maximizando a evidência

- Lembrando a diferença entre ML e ML-II:
 - Máxima verossimilhança (ML):

$$\hat{\mathbf{w}} = \arg \max \log p(\mathcal{D}|\mathbf{w}).$$

- Máxima verossimilhança do tipo II (ML-II):

$$\hat{\phi} = \arg \max \log \int p(\mathcal{D}|\mathbf{w}, \phi) p(\mathbf{w}|\phi) d\mathbf{w}.$$

- Encontrar os hiperparâmetros $\hat{\phi}$ via ML-II pode ser feito por:
 - Uso dos gradientes da evidência.
 - Algoritmo Expectation-Maximization (EM).
 - Computar a evidência para uma grid de candidatos.

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

Tópicos adicionais

- Regressão Bayesiana com variância σ^2 do ruído desconhecida com priori Gamma inversa $\text{IG}(\sigma^2|a_0, b_0)$:

$$p(\mathbf{y}|\Phi, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \sigma^2\mathbf{I}), \quad \Phi = \phi(\mathbf{X}),$$
$$p(\mathbf{w}, \sigma^2|\mathbf{m}_0, \mathbf{S}_0, a_0, b_0) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)\text{IG}(\sigma^2|a_0, b_0).$$

- A posteriori $p(\mathbf{w}, \sigma^2|\mathcal{D})$ é analítica (Murphy, pgs. 235 e 236), sendo as marginais $p(\mathbf{w}|\mathcal{D})$ e $p(\sigma^2|\mathcal{D})$ respectivamente uma t de Student e uma Gamma inversa.
- A preditiva $p(\mathbf{y}_*|\mathbf{X}_*)$ também é analítica (Murphy, pg. 236), sendo uma t de Student.
- Modelos não-conjugados.

Agenda

- 1 Regressão linear iterativa
 - Regressão linear simples
 - Regressão linear múltipla
- 2 Regressão linear analítica
- 3 Regressão Linear Bayesiana
- 4 Comparação Bayesiana de modelos
- 5 Tópicos adicionais
- 6 Referências

Referências bibliográficas

- **Cap. 9** - DEISENROTH, M. *et al.* **Mathematics for machine learning**. 2019.
- **Caps. 1 e 7** - MURPHY, Kevin P. **Machine learning: a probabilistic perspective**, 2012.
- **Caps. 2 e 3** - HAYKIN, Simon. **Neural Networks and Learning Machines**, 3ed., 2010.
- **Cap. 3** - BISHOP, Christopher M. **Pattern recognition and machine learning**, 2006.