#### Lecture 23 Bayesian regression methods

Bayes' rule; weight-based view and linear model

Prof. David A. Kofke CE 500 – Modeling Potential-Energy Surfaces Department of Chemical & Biological Engineering University at Buffalo



© 2024 David Kofke

## What color is the candy that I will select from this urn?

- It won't be white
- It is one of these colors black, red, green, blue, yellow, brown, orange, purple, gray
- The candy is not licorice
- The candy is M&Ms
- The color distribution of M&Ms is 24% blue, 20% orange, 16% green, 14% yellow, 13% red, and 13% brown
- What principle are you using to make your prediction?

# Probability (like temperature ) is a more subtle concept than you might realize

- Two interpretations are in common use
- Frequentist
  - Probability describes the number of times each outcome would occur if an infinite number of samples were taken on a population
  - It can change only by changing the population
- Bayesian

4

- Probability quantifies how much information you have about an outcome
- It can change as more information becomes available

J

- Here is a space of outcomes of an event:
  - Select a number uniformly at random from 1 to 100

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34
	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51
	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68
	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85
	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	D	
C	•																

- Events A and B occur when numbers are selected in the indicated regions, respectively
- What are P(A)? P(B)?

E																		
		7																
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	<b>)</b> 17
	-	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34
		35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51
	5	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68
	(	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85
	٤	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	)	
	S																	

- Events A and B occur when numbers are selected in the indicated regions, respectively
- What are P(A) = 34/100 = 0.34P(B) = 12/100 = 0.12

<b>A</b> 1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	<b>3</b> 17
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34
35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51
52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68
69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85
86	87	88	89	90	91	92	93	94	95	96	97	98	99	100	)	
S																

• What is the probability of B, if we know that A occurred? Probability of B, given A P(B|A) = ?

<b>A</b> 1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	<b>3</b> 17
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34

• What is the probability of B, if we know that A occurred? Probability of B, given A P(B|A) = 8/34 = 0.24

<b>A</b> 1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	<b>3</b> 17
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34

• What is the probability of B, if we know that A occurred? Probability of B, given A P(B|A) = 8/34 = 0.24 $= P(A \cap B)/P(A)$ =(8/100)/(34/100)

<b>A</b> 1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	<b>3</b> 17
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34

• This is the general formula for the *conditional probability* 

A 1	2	2	Λ	5	6	7	8	٩	10	11	12	12	1/	15	16	<b>3</b> 17
1	۷	5	т	5	U	,	0	5	10	11	12	13	74	13	10	1/
18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34

$$P(\mathrm{B}|\mathrm{A}) = rac{P(\mathrm{B}\cap\mathrm{A})}{P(\mathrm{A})}$$

11

- It works both ways
- What is P(A|B)?

A	14	15	16	<b>3</b> 17
	31	32	33	34
	48	49	50	51

- It works both ways • What is P(A|B) = 8/12 = 0.67  $P(A \cap B)$ ? P(B)? A A B 14 15 16 17 31 32 33 34 48 49 50 51
  - $P(A \cap B)/P(B)$ ?

•	It works both ways	Α		B	
•	What is		14 15	16 1	
	P(A B) = 8/12 = 0.67		31 32	33 3	34
	$P(A \cap B) = 8/100$		48 49	50 5	51
	P(B) = 12/100				
	$P(A \cap B)/P(B) = (8/100)$	(12/100) = 8/12			

 We have expressions for the two conditional probabilities
 P(B|A) = 
 P(B ∩ A)
 P(A|B) = 
 P(A ∩ B)
 P(B)

 Or

 $P(B \cap A) = P(B|A)P(A)$   $P(A \cap B) = P(A|B)P(B)$ 

• Of course,  $(A \cap B)$  and  $(B \cap A)$  are the same event, so

 $P(\mathrm{B}|\mathrm{A})P(\mathrm{A})=P(\mathrm{A}|\mathrm{B})P(\mathrm{B})$ 

– This is *Bayes' rule* 

# Bayes' rule provides a framework to update probabilities when given new information

#### $P(\mathbf{B}|\mathbf{A})P(\mathbf{A}) = P(\mathbf{A}|\mathbf{B})P(\mathbf{B})$

• It's use is best understood when written in this form



• The posterior is the updated probability in light of new information

## The various parts of Bayes' formula can be explained in words

- In our applications, event A is a model parameter(s) w or outcome y having a particular value, and event B is data
- P(w) this is the *prior probability*, the probability of the parameter value (or outcome y) in the absence of data
   Could be a uniform distribution over an expected range
- P(w|data) this is the *posterior probability*, informed by data
- P(data|w) this is the *likelihood* that we'd observe the given data, given that the model parameter has the specified value
- P(data) this is the *marginal probability* of the data, considering all possible parameter values

# The marginal probability is usually obtained by integration, or normalization

$$P(\mathbf{w}| ext{data}) = P(\mathbf{w}) imes rac{P( ext{data}|\mathbf{w})}{P( ext{data})}$$



• Consider all mutually exclusive B events

$$egin{aligned} P(\mathrm{A}) &= P(\mathrm{A} \cap \mathrm{B}_1) + \dots + P(\mathrm{A} \cap \mathrm{B}_n) \ &= P(\mathrm{A}|\mathrm{B}_1)P(\mathrm{B}_1) + \dots + P(\mathrm{A}|\mathrm{B}_n)P(\mathrm{B}_n) \end{aligned}$$

• Or more generally for our case  $P(\text{data}) = \int P(\text{data}|\mathbf{w})P(\mathbf{w})d\theta$ The marging depend on

18

The marginal probability does not depend on **w** 

# Regression is a process of induction: specific $\rightarrow$ general

- Given a finite set of training data
- We want to infer a function or model *f*(**x**) that provides accurate outputs for all possible input values
- The data are insufficient to do this perfectly, i.e., to solve for f
- Hence, we impose conditions, or assumptions, regarding the form of *f*

- Otherwise, any function that goes through the points would be valid

#### Two approaches to restrict the model are in common use

- Restrict the class of functions that are considered
  - e.g., linear functions of the variables in **x**
  - Requires a decision on the richness of the class of functions considered
    - Target may not be well described by the selected functions
    - Adding more functions risks overfitting
- Assign a *prior probability* to all trial functions
  - Higher probability assigned to more likely functions, e.g., those that are smooth versus those that are erratic
  - This admits and infinite range of functions; how to account for all of them? (Gaussian processes does the trick)

#### Rather than going straight to functions, a weightspace view can ease into the Bayesian approach

• Standard linear model



• Gaussian noise is added to linear model to introduce deviations

 $f(\mathbf{x}) = \mathbf{x}^ op \mathbf{w} \qquad y = f(\mathbf{x}) + arepsilon \qquad \mathbf{y} = X\mathbf{w} + oldsymbol{arepsilon} \qquad arepsilon \sim \mathcal{N}(0, \sigma_n^2)$ 

# We start by writing the likelihood of the observed data, for the given inputs and weights

• Assuming independent data:

Deviation from model follows Gaussian pdf

$$egin{split} p(\mathbf{y} \mid X, \mathbf{w}) &= \prod_{i=1}^n p(y_i \mid \mathbf{x}_i, \mathbf{w}) = \prod_{i=1}^n rac{1}{\sqrt{2\pi}\sigma_n} \expigg(-rac{\left(y_i - \mathbf{x}_i^ op \mathbf{w}
ight)^2}{2\sigma_n^2}igg) \ &= \left(2\pi\sigma_n^2
ight)^{-n/2} \expigg(-rac{1}{2\sigma_n^2}|\mathbf{y}-X\mathbf{w}|^2igg) = \mathcal{N}ig(X\mathbf{w}, \sigma_n^2Iig) \end{split}$$

- We need to specify a prior probability for the weights
  - Choose a zero-mean Gaussian with covariance matrix  $\Sigma_p$

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_{\mathrm{p}})$$

# The Bayesian linear model specifies the weights via a posterior distribution, given the data

$$p(\mathbf{w} \mid \mathbf{y}, X) = p(\mathbf{w}) rac{p(\mathbf{y} \mid X, \mathbf{w})}{p(\mathbf{y} \mid X)}$$
 .



- The marginal probability is independent of **w**
- Obtainable as normalization constant:  $p(\mathbf{y} \mid X) = \int p(\mathbf{y} \mid X, \mathbf{w}) p(\mathbf{w}) d\mathbf{w}$
- The result can be expressed as a multivariate Gaussian

$$p(\mathbf{w} \mid \mathbf{y}, X) \sim \mathcal{N}(ar{\mathbf{w}}, A^{-1})$$

- Mean:  $\mathbf{\bar{w}} = \frac{1}{\sigma_n^2} A^{-1} X^{\top} \mathbf{y}$
- Covariance matrix:  $A^{-1}$ ,  $A = \sigma_n^{-2} X^{\top} X + \Sigma_p^{-1}$

#### The recommended weights can be extracted in several ways

- Use the mean (expected value), or the mode (most likely value)
  These are the same for the Normal distribution
- We obtain:  $\mathbf{\bar{w}} = \left(X^{\top}X + \sigma_n^2\Sigma_{\mathrm{p}}^{-1}\right)^{-1}X\mathbf{y}$
- Compare to result for ridge regression (Lecture 19)

$$\widehat{\mathbf{w}} = (X^ op X + \lambda I)^{-1} X^ op \mathbf{y}$$

• The prior on the weights imposes a regularization penalty!

## **Everything is given as a probability distribution**



Rasmussen & Williams

#### To make predictions in the Bayesian scheme, we can average over the entire distribution of weights

• This again results in a Gaussian for the predicted values

$$p(f_* \mid \mathbf{x}_*, X, \mathbf{y}) = \int p(f_* \mid \mathbf{x}_*, \mathbf{w}) p(\mathbf{w} \mid X, \mathbf{y}) d\mathbf{w}$$

$$predicted value = \mathcal{N}(\underbrace{\sigma_n^{-2} \mathbf{x}_*^\top A^{-1} X^\top \mathbf{y}}_{\text{vector}}, \underbrace{\mathbf{x}_*^\top A^{-1} \mathbf{x}_*}_{\text{vector}})$$

$$= \mathcal{N}(\underbrace{\sigma_n^{-2} \mathbf{x}_*^\top A^{-1} X^\top \mathbf{y}}_{\text{mean}}, \underbrace{\mathbf{x}_*^\top A^{-1} \mathbf{x}_*}_{\text{variance}})$$

$$A = \sigma_n^{-2} X^\top X + \Sigma_p^{-1}$$

$$\sigma_n^{-2} A^{-1} = (X^\top X)^{-1}$$
and the mean  $f_*(\mathbf{x}_*)$  is  $\mathbf{x}_*^\top (X^\top X)^{-1} X^\top \mathbf{y} = \mathbf{x}_*^\top \mathbf{\bar{w}}$ 

# Features (a.k.a. basis functions) are a way to add versatility to the model formulation

- Expand the *p*-dimensional **x** vector into an *N*-dimensional space of variables φ(**x**)
- E.g., instead of a straight-line fit to a single x variable, use a polynomial fit  $\begin{pmatrix} 1 \\ \end{pmatrix}$

$$\mathbf{x} = (x) \qquad oldsymbol{\phi}(\mathbf{x}) = egin{pmatrix} x \ x^2 \ dots \$$

$$f(\mathbf{x}) = \boldsymbol{\phi}(\mathbf{x})^{\top} \mathbf{w}$$

#### Analysis using features $\phi$ proceeds exactly as for case working directly with input vector x

• Replace *X* with feature matrix

$$\Phi(X) = egin{pmatrix} \phi_1^{(1)}(\mathbf{x}) & \dots & \phi_N^{(1)}(\mathbf{x}) \ dots & \ddots & dots \ \phi_1^{(n)}(\mathbf{x}) & \dots & \phi_N^{(n)}(\mathbf{x}) \end{pmatrix}$$

 $n \times N$  matrix

• The predictive distribution is

$$egin{aligned} &f_* \mid \mathbf{x}_*, X, \mathbf{y} \sim N(\sigma_n^{-2} oldsymbol{\phi}_*^ op A^{-1} \Phi \mathbf{y}, oldsymbol{\phi}_*^ op A^{-1} oldsymbol{\phi}_*) \ &oldsymbol{\phi}_* \equiv oldsymbol{\phi}(\mathbf{x}_*) & A = \sigma_n^{-2} \Phi^ op \Phi + \Sigma_\mathrm{p}^{-2} \end{aligned}$$

## The result may be reformulated to introduce a *kernel* in lieu of the features

- Inversion of  $N \times N$  matrix A may be inconvenient
- The distribution can be rewritten as

$$egin{aligned} f_* \mid \mathbf{x}_*, X, \mathbf{y} &\sim \mathcal{N}(oldsymbol{\phi}_*^ op \Sigma_p \Phi^ op ig(K + \sigma_n^2 Iig)^{-1} \mathbf{y} \ oldsymbol{\phi}_*^ op \Sigma_p oldsymbol{\phi}_* - oldsymbol{\phi}_*^ op \Sigma_p \Phi^ op ig(K + \sigma_n^2 Iig)^{-1} \Phi \Sigma_p oldsymbol{\phi}_* ig) \ & oldsymbol{-} \end{aligned}$$

- Where the kernel matrix is defined  $K = \Phi \Sigma_p \Phi^{+}$
- This instead requires inversion of an  $n \times n$  matrix
- Desirable if feature space is larger than amount of data N > n

## The result may be reformulated to introduce a *kernel* in lieu of the features

- Inversion of  $N \times N$  matrix A may be inconvenient
- The distribution can be rewritten as

$$egin{aligned} f_* \mid \mathbf{x}_*, X, \mathbf{y} &\sim \mathcal{N}igg( oldsymbol{\phi}_*^ op \Sigma_p \Phi^ op igl( K + \sigma_n^2 I igr)^{-1} \mathbf{y} \ & oldsymbol{\phi}_*^ op \Sigma_p oldsymbol{\phi}_* igr] - oldsymbol{\phi}_*^ op \Sigma_p \Phi^ op igl( K + \sigma_n^2 I igr)^{-1} oldsymbol{\Phi} \Sigma_p oldsymbol{\phi}_* igr) \end{aligned}$$

- Note that the features always enter in a specific way
- Define covariance function, or *kernel*  $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x}) \Sigma_p \boldsymbol{\phi}(\mathbf{x}')$
- Express as a dot product:  $k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\psi}(\mathbf{x}) \cdot \boldsymbol{\psi}(\mathbf{x}'), \quad \boldsymbol{\psi}(\mathbf{x}) \equiv \Sigma_p^{1/2} \boldsymbol{\phi}(\mathbf{x})$

# The *kernel trick* replaces inner products in x with the kernel operation, k(x,x')

- This allows model to bypass feature vectors entirely in lieu of the kernel
- The kernel then becomes the central focus of model development

#### **Suggested Reading/Viewing**

- C. E. Rasmussen & C. K. I. Williams, *Gaussian Processes for Machine Learning*, the MIT Press, 2006. Chapters 1 and 2.
  - www.GaussianProcess.org/gpml
- *Quantum Chemistry in the Age of Machine Learning*, edited by P. O. Dral
  - Chapter 9. Kernel Methods, Max Pinheiro Jr. and Pavlo O. Dral
  - Chapter 10. Bayesian Inference, Wei Liang and Hongsheng Dai
  - Posted on UBLearns