

Course: Introduction to Machine Learning

Assignment 4

Important NOTE: In order to get full credit, for every question, you need to provide the details of your work on how to get to a solution or the end of a proof.

Instructor: Tan Bui-Thanh

TA: Cole Nockolds

Due day: 11:59 pm, 29 October, Friday

All submission MUST be in pdf format, except codes. All code outputs should be reported in pdf format.

Question 1 (30 points) Probabilistic linear regression

Recall: The multivariate Gaussian distribution has the form $\mathcal{N}(\mu, \Sigma)$, $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$

$$\mathcal{N}(\boldsymbol{\mu}, \Sigma) = (2\pi)^{-\frac{n}{2}} \det(\Sigma)^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \propto \underbrace{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)}_{(\propto \text{ is to ignore the normalization constant})}$$

(\propto is to ignore the normalization constant).

To find the mean $\boldsymbol{\mu}$, we take first derivative of $-\log(F)$ and solve for \mathbf{x} , i.e.,

$$\frac{d}{d\mathbf{x}}[-\log(F)] = \mathbf{0} \implies \mathbf{x} = \boldsymbol{\mu}$$

To find the covariance matrix, Σ , we take the second derivative of $-\log(F)$, which returns Σ^{-1} , then invert the matrix to obtain Σ .

$$\Sigma = \left[\frac{d^2}{d\mathbf{x}^2} [-\log(F)] \right]^{-1}$$

By these procedures, we can find the mean $\boldsymbol{\mu}$ and the covariance matrix Σ of original multivariate Gaussian distribution. In the case that the distribution is not directly the form of Gaussian distribution, for example,

$$F = \exp\left(-\frac{1}{2}(\mathbf{x}^T (A^T A) \mathbf{x} + \mathbf{x}^T \mathbf{A})\right)$$

following the same procedure, you would find the mean is $\mathbf{0}$, and the covariance matrix is $(A^T A + I)^{-1}$.

Now, generate 6 data samples for the problem using the polynomial function (use the below given code for create $x^{(i)}$)

$$f(x) = .5 - x - .5x^2 - 2x^3 + 5x^4.$$

Then add noise: $y^{(i)} = f(x^{(i)}) + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \delta^2)$, and noise level $\delta = 0.2$. In this problem, we assume that the model is 2nd order, i.e., PolynomialFeatures(degree=2) in scikit-learn. \

1. Using a Bayesian approach, find, analytically, the mean and covariance of the posterior distribution $P(y|x, \mathcal{D})$, where, as in lecture 5b, we have

$$P(y|x, \mathcal{D}) = \int_{\theta} P(y | x, \boldsymbol{\theta}) P(\boldsymbol{\theta} | \mathcal{D}) d\boldsymbol{\theta}.$$

where

- $P(y | x, \boldsymbol{\theta})$ is the likelihood distribution (chosen to be the set of a Gaussian distributions as in lecture 5b)

- $P(\boldsymbol{\theta} | \mathcal{D})$ is the posterior distribution given the training data, as in lecture 5a, we have

$$P(\boldsymbol{\theta} | \mathcal{D}) = \frac{P(\mathcal{D} | \boldsymbol{\theta})P(\boldsymbol{\theta})}{P(\mathcal{D})}$$

in which $P(\boldsymbol{\theta})$ is the Gaussian prior with mean zero and variance σ^2 .

Hints:

- Show that the distribution $P(\boldsymbol{\theta} | \mathcal{D}) \sim \mathcal{N}(\boldsymbol{\theta}_{\text{MAP}}, \Sigma_{\boldsymbol{\theta}})$ has the form of Gaussian distribution.
 - Apply the multivariate Gaussian distribution formula to find $\boldsymbol{\theta}_{\text{MAP}}, \Sigma_{\boldsymbol{\theta}}$.
 - Find the distribution $P(y|x, \mathcal{D})$ by marginalizing the $\boldsymbol{\theta}$ random variable.
2. Plot the uncertainty band with the range of $x \in [0, 1]$.
 3. Discuss the benefits of a Bayesian method as opposed to maximum likelihood.

1.1 - derivation of $P(y|x, D)$

Anushi R. Sadam

Professor Tan Bui Thanh

COE 379L: Introduction to Machine Learning & Data Science

29 October 2024

Homework 4

- Find the mean, covariance of $P(y|x, D)$
w/ Bayesian approach.

$$P(y|x, D) = \int_{\theta} P(y|x, \theta) \cdot P(\theta|D) d\theta$$

$$P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$$

we assume the prior $P(\theta) = N(0, \sigma^2 I_d)$
(mean 0, variance σ^2), since we have
no info about the prior in our test problem.
 d is the dimension of θ ; $\theta \in \mathbb{R}^d$.

$P(D|\theta)$ is the likelihood based on the dataset

$$D = \{(x^i, y^i)\}_{i=1}^n \sim \text{number of samples} = |D|$$

$$P(D|\theta) = \prod_{i=1}^n P(x^i, y^i | \theta)$$

$$\text{we know } P(x, y | \theta) = P(y|x, \theta) \cdot P(x|\theta)$$

Since $x \perp \theta$ (x is independent of θ), $P(x|\theta) = P(x)$. $P(x)$ is an uniform distribution since we must assume our data collection evenly covers the system.

We assume a gaussian conditional distribution:

2a

$$P(y|x_i, \theta) = N(\theta^T x_i, \sigma^2)$$

feature vector

or: (n data samples)

$$P(y|x, \theta) = \prod_{i=1}^n N(\theta^T x_i, \sigma^2)$$

assume same noise level
for all measurements.

$$\text{Thus, } P(x, y | \theta) \propto P(y|x, \theta)$$

$$\propto \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y_i - \theta^T x_i)^2\right)$$

$$P(x, y | \theta) = P(y|x, \theta) \cdot P(x|\theta)$$

Since x is independent of θ , $P(x|\theta)$ is $P(x)$.

So, since $P(x)$ is a uniform distribution,

$$P(x, y | \theta) \propto P(y|x, \theta)$$

normal distribution assumed by our model.

$$P(D|\theta) = P(x, y | \theta) \propto \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y_i - \theta^T x_i)^2\right)$$

We assume $P(\theta)$ is a gaussian distribution w/ mean $\vec{\theta}$, variance of each element λ^2 :

(2b)

Now, we assume a prior of mean 0, variance $\lambda^2 = 10$
 (along all dimensions)

in general:

$$P(\vec{\theta}) = \frac{1}{\sqrt{(2\pi)^{|D|} |\Sigma|^{\frac{1}{2}}}} \exp\left(-\frac{1}{2} (\vec{\theta} - \vec{\mu})^\top \Sigma^{-1} (\vec{\theta} - \vec{\mu})\right)$$

where $|D|$ is the number of dimensions for $\vec{\theta}$,
 Σ is the covariance matrix,
 μ is the mean of $P(\vec{\theta})$.

We know $\mu = \vec{0}$. Now to calculate Σ :

$$\Sigma_{ij} = \text{cov}(\theta_i, \theta_j) \quad \forall i, j \in [1, |D|]$$

$$= \mathbb{E}[(\theta_i - \mathbb{E}[\theta_i])(\theta_j - \mathbb{E}[\theta_j])]$$

$$\mathbb{E}[\theta_i] = \mathbb{E}[\theta_j] = 0 \quad (\text{mean is } \vec{0})$$

$$\Rightarrow \Sigma_{ij} = \mathbb{E}[\theta_i \times \theta_j]$$

i) $\forall i=j$:

$$\Sigma_{ii} = \mathbb{E}[\theta_i^2]$$

$$\text{we know } \text{Var}(\theta_i) = \mathbb{E}[\theta_i^2] - \underbrace{(\mathbb{E}[\theta_i])^2}_{0} = \lambda^2$$

$$\Rightarrow \text{Var}(\theta_i) = \mathbb{E}[\theta_i^2] = \lambda^2$$

for each θ_i , variance
is a constant, λ^2

Thus, $\Sigma_{ii} = \lambda^2$

2) $i \neq j$

$$\Sigma_{ij} = E[\theta_i \theta_j]$$

we assume independent variables for the coefficients θ_i random

$$\text{so, } \Sigma_{ij} = E[\theta_i] \cdot E[\theta_j]$$

$$= \mu_i \cdot \mu_j$$

mean \Rightarrow mean of θ_j
of θ_i distribution distribution

Since we assume random normal distributions of mean 0 ($\mu_i = 0, \mu_j = 0$).

$$\Sigma_{ij} = 0 \cdot 0 = 0.$$

$$\text{Thus, } \Sigma = \begin{bmatrix} \lambda^2 & 0 & \dots & 0 \\ 0 & \lambda^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda^2 \end{bmatrix} = \lambda^2 I_{|D|}$$

identity of dimension
 $|D| \times |D|$

$$\det(\Sigma) = (\lambda^2)^{|D|}$$

$$\Sigma^{-1} = (\lambda^2 I_{|D|})^{-1} = \frac{1}{\lambda^2} I_{|D|}$$

Note $P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$

3

$P(D)$ is simply a normalization value used so that the total probability $\int_D P(\theta|D) d\theta = 1$.

($\Rightarrow P(\theta|D)$ is a bona fide probability distribution)

Then,

$$P(\theta|D) \propto P(D|\theta) \cdot P(\theta)$$

$$= \prod_{i=1}^n \exp\left(-\frac{1}{2\sigma^2} (y^i - \theta^T x^i)^2\right) \cdot$$

$$\exp\left(-\frac{1}{2\lambda^2} \|\theta\|_2^2\right)$$

$$= \exp\left(-\frac{1}{2} \sum_{i=1}^n \left[\frac{1}{\sigma^2} (y^i - \theta^T x^i)^2 + \frac{1}{\lambda^2} \|\theta\|_2^2 \right]\right)$$

$$P(\theta|D)$$

we seek to maximize \propto feature matrix

$$P(\theta|D) \propto \exp\left(-\frac{1}{2} \left[(y - x^\top \theta)^2 \cdot \frac{1}{s^2} + \frac{1}{\lambda^2} \theta^\top \theta \right]\right)$$

$$\propto \exp\left(-\frac{1}{2} \left[\frac{1}{s^2} (y^\top y - 2y^\top x \theta + \theta^\top x^\top x \theta) + \frac{1}{\lambda^2} (\theta^\top \theta) \right]\right)$$

$$\propto \exp\left(-\frac{1}{2} \underbrace{\left(\theta^\top \left(\frac{x^\top x}{s^2} + \frac{I}{\lambda^2} \right) \theta - 2 \frac{y^\top x}{s^2} \theta + \frac{y^\top y}{s^2} \right)}_{A^\top A - b^\top b + c}\right)$$

$$\text{let } A = \begin{bmatrix} x^\top x & I \\ s^2 & \lambda^2 \end{bmatrix} \quad \& \quad b = \frac{x^\top y}{s^2} \quad \& \quad c = \frac{y^\top y}{s^2}$$

then:

$$\begin{aligned} -\frac{1}{2} (\theta^\top A \theta - 2b^\top \theta + c) &= -\frac{1}{2} (\theta^\top A \theta - 2b^\top \theta + b^\top A^{-1} b \\ &\quad - b^\top A^{-1} b + c) \\ &= -\frac{1}{2} ((\theta - A^{-1} b)^\top A (\theta - A^{-1} b) - b^\top A^{-1} b + c) \end{aligned}$$

FA side:

$$\theta^\top A \theta - 2b^\top \theta = (\theta^\top A \theta - 2b^\top \theta + b^\top A^{-1} b) - b^\top A^{-1} b - 2b^\top \theta$$

$$\theta^\top A \theta - 2b^\top \theta + b^\top A^{-1} b = (\theta - A^{-1} b)^\top A (\theta - A^{-1} b)$$

since:

LHS expansion:

$$\begin{aligned} (\theta^\top - b^\top) A \theta - (\theta^\top - b^\top) A \cdot A^{-1} b - (A^{-1} b)^\top A \theta \\ + (A^{-1} b)^\top A (A^{-1} b) \end{aligned}$$

$$(A^{-1}b)^T A = b^T (A^{-1})^T A.$$

5

$$\text{Note since } I = A^{-1}A = AA^{-1}; \quad I^T = (A^{-1}A)^T = (AA^{-1})^T$$

$$I = I^T = (AA^{-1})^T = (A^{-1})^T A^T$$

Note that A is a symmetric matrix $A = mx^T x + nI$
where m, n are constants

$$A^T = ((mx^T x) + nI)^T = (nI)^T + (mx^T x)^T$$

$$= nI^T + m(x^T x)^T = nI + m x^T (x^T)^T$$
$$= nI + mx^T x = mx^T x + nI = A.$$

Thus, A is symmetric.

$$\text{so, } I = (A^{-1})^T A^T = (A^{-1})^T A.$$

$$\text{Then, } (A^{-1}b)^T A = b^T I = b^T.$$

Continuing the LHS expansion:

$$\theta^T A \theta - \theta^T b - b^T \theta + \underbrace{(b^T (A^{-1})^T A A^{-1} b)}_I$$

$$= \theta^T A \theta - 2b^T \theta + b^T (A^{-1}) b \quad \exists$$

$$\text{Thus, } P(\theta | D) \propto \exp\left[-\frac{1}{2}\left((\theta - A^{-1}b)^T A (\theta - A^{-1}b) - b^T A^{-1} b + c\right)\right]$$

continuing w/ our exponential factoring:

$$\text{let } F = (b^T A^{-1} b + c) \frac{1}{2}$$

then in the exponent, we have

b

$$-\frac{1}{2} (\theta - \theta^{-1} b)^T A (\theta - \theta^{-1} b) + F.$$

so:

$$\begin{aligned} P(\theta | D) &\propto \exp\left(-\frac{1}{2} (\theta - \theta^{-1} b)^T A (\theta - \theta^{-1} b) + F\right) \\ &\propto \exp\left(-\frac{1}{2} \left((\theta - \theta^{-1} b)^T A (\theta - \theta^{-1} b) \right)\right) \cdot e^F \\ &\propto \exp\left(-\frac{1}{2} (\theta - \theta^{-1} b)^T A (\theta - \theta^{-1} b)\right) \end{aligned}$$

multivariate Gaussian expression 

w/ mean $\theta^* = \theta^{-1} b$, $\Sigma_0 = \theta^{-1}$.

$$\text{Thus, } P(\theta | D) \propto \exp\left(-\frac{1}{2} (\theta - \theta^*)^T \Sigma_0^{-1} (\theta - \theta^*)\right)$$

$$\text{where } \theta^* = \left[\frac{\mathbf{x}^T \mathbf{x}}{\sigma^2} + \frac{\mathbf{I}}{\lambda^2} \right]^{-1} \frac{\mathbf{x}^T \mathbf{y}}{\sigma^2} = \left[\mathbf{x}^T \mathbf{x} + \frac{\sigma^2}{\lambda^2} \mathbf{I} \right]^{-1} \mathbf{x}^T \mathbf{y}$$

$$\Sigma_0 = \left[\frac{\mathbf{x}^T \mathbf{x}}{\sigma^2} + \frac{\mathbf{I}}{\lambda^2} \right]^{-1} = \sigma^2 \left[\mathbf{x}^T \mathbf{x} + \frac{\sigma^2}{\lambda^2} \mathbf{I} \right]^{-1}$$

① Note about scalar inverse pg 7

we recognize θ^* as a least squares solution w/ regularization

$$\frac{\sigma^2}{\lambda^2}$$

let $\hat{y} = \mathbf{x}\theta$, then we have a least squares problem:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} (y - \mathbf{x}\theta)^T (y - \mathbf{x}\theta) + \frac{\gamma}{2} \|\theta\|^2, \gamma \text{ is a regularizer}$$

The solution to the L_2 regularized solution is

loss:

$$J(\theta) = \|y - X\theta\|^2 + \frac{\gamma}{2} \|\theta\|^2.$$

$$\nabla J(\theta) = X^T X \theta - X^T y + \gamma \theta = 0.$$

$$\underbrace{\theta(X^T X + \gamma I)}_{\text{positive definite}} = X^T y$$

$$\text{positive definite; } \theta_{\text{MAP}}^* = (X^T X + \gamma I)^{-1} X^T y$$

$$\text{Thus, for } P(\theta | D), \theta^* = \left(X^T X + \frac{\gamma^2}{\lambda^2} I \right)^{-1} X^T y$$

which is simply the MAP solution for θ^* via least squares, w/ regularizer $\gamma = \frac{\gamma^2}{\lambda^2}$.

Note that the mean of $P(\theta | D)$ is θ^* , since

$$P(\theta | D) \propto \exp\left(-\frac{1}{2} (\theta - \theta^*)^T \Sigma_0 (\theta - \theta^*)\right)$$

has the same form as the gaussian w/
mean $\mu = \theta^*$ & covariance Σ_0 .

$$\Sigma_0, \text{ thus, is: } \boxed{\Sigma_0 = \left[\frac{X^T X}{\gamma^2} + \frac{I}{\lambda^2} \right]^{-1} = \frac{\gamma^2}{\lambda^2} \left[X^T X + \frac{\gamma^2}{\lambda^2} I \right]^{-1}}$$

① $(CA)^{-1} = C^{-1}A^{-1}$ if $C \in \mathbb{R}^{d \times d}$, $A \in \mathbb{R}^{d \times d}$ def in \mathbb{R}
 $\downarrow A$ must be invertible

Proof:

Suppose CA has inverse B :

$$(CA)B = I; C(CAB) = I; CAB = \underset{C}{\cancel{C}} I$$

$$A^{-1}AB = \underset{C}{\cancel{C}} A^{-1} I; IB = \underset{C}{\cancel{C}} A^{-1}$$

$$\therefore B = \underset{C}{\cancel{C}} A^{-1} = C^{-1}A^{-1}$$

Note that $P(\theta|D) = N(\theta^*_{MAP}, \Sigma_0)$; 8
 we can change from \propto to = since
 we know $P(\theta|D)$ must be a bona fide
 probability distribution, so, $\int_0 P(\theta|D) d\theta = 1$.

If $P(\theta|D) = c \cdot N(\theta^*_{MAP}, \Sigma_0)$ for some
 constant (scalar) c , then,

$$\int_0 P(\theta|D) d\theta = \int_{-\infty}^{\infty} c \cdot \frac{\exp(-\frac{1}{2}(\theta - \theta^*_{MAP})^T \Sigma_0^{-1} (\theta - \theta^*_{MAP}))}{(2\pi)^{d/2} |\Sigma_0|^{1/2}} d\theta$$

where d is the dimension of the data

$$\text{we know, then } \int_0 P(\theta|D) d\theta = c \cdot \int_{-\infty}^{\infty} N(\theta^*_{MAP}, \Sigma_0) d\theta.$$

per the definition of a normal distribution,
 we know

$$\int_{-\infty}^{\infty} N(\theta^*_{MAP}, \Sigma_0) d\theta = 1$$

(area under normal distribution = 1).

$$\text{Then, } \int_0 P(\theta|D) d\theta = c \cdot 1. \text{ so, } c = 1.$$

This means $P(\theta|D) = N(\theta^*_{MAP}, \Sigma_0)$.

The coefficient $\frac{1}{(2\pi)^{d/2} |\Sigma_0|^{1/2}}$ comes

from the various constants we ignored in
 our proportion analysis done previously.

For example, one of the constants being
 ignored is $P(D)$ when we solve for $P(\theta|D)$.

This will help allow $P(\theta|D)$ to be normalized.

1. cont.

Do 3. Now, to find $P(y|x, D)$:

$$\begin{aligned} P(y|x, D) &= \int P(y|x, \theta) P(\theta|D) d\theta \\ &= \int_{\theta} \prod_{i=1}^n N(\theta^T x_i, \delta) \cdot N(\theta_{MAP}, \Sigma_0) d\theta \end{aligned}$$

note that

$$\prod_{i=1}^n N(\theta^T x_i, \delta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\delta^2} \exp\left(-\frac{1}{2\delta^2}(y_i - \theta^T x_i)^2\right)$$

$$= \frac{1}{(2\pi)^{\frac{n}{2}}\delta^n} \exp\left(-\frac{1}{2\delta^2}(y - x\theta)^T(y - x\theta)\right)$$

where X is the feature Matrix

$$\text{so, } P(y|x, \theta) \cdot P(\theta|D) =$$

$$\begin{aligned} &\left(\frac{1}{(2\pi)^{\frac{n}{2}}\delta^n} \cdot \frac{1}{(2\pi)^{\frac{n}{2}}|\Sigma_0|} \right) \underbrace{\exp\left(-\frac{1}{2\delta^2}(y - x\theta)^T(y - x\theta)\right)}_{M} \\ &\quad + \underbrace{-\frac{1}{2}(\theta - \theta_{MAP})^T \Sigma_0^{-1} (\theta - \theta_{MAP})}_{\text{constant}} \end{aligned}$$

let M be the exponent term. M we expand M to get:

$$\begin{aligned} &\frac{-1}{2\delta^2} (y^T y - 2y^T x\theta + \theta^T x^T x\theta) - \frac{1}{2} (\theta^T \Sigma_0^{-1} \theta - 2\theta^T \Sigma_0^{-1} \theta_{MAP} \\ &\quad + \theta_{MAP}^T \Sigma_0^{-1} \theta_{MAP}) \end{aligned}$$

Grouping quadratic terms:

10

$$M = -\frac{1}{2} \left(\theta^T \left(\frac{x^T x}{s^2} + \Sigma_0^{-1} \right) \theta \right)$$

$$-2\theta^T \left(\frac{x^T y}{s^2} + \Sigma_0^{-1} \theta_{MAP} \right)$$

$$\underbrace{-\frac{1}{2s^2} y^T y - \frac{1}{2} \theta_{MAP}^T \Sigma_0^{-1} \theta_{MAP}}_{\text{constant } C \text{ wrt } \theta}$$

C is a constant so, we can factor out e^C as a constant from the integral of $P(y|x, D)$.

$$\text{Thus, define } Q = -\frac{1}{2} (\theta^T A \theta - 2\theta^T b)$$

be the exponent in the integral,

$$\text{where } A = \frac{x^T x}{s^2} + \Sigma_0^{-1} \quad \& \quad b = \frac{x^T y}{s^2} + \Sigma_0^{-1} \theta_{MAP}.$$

As before, we can complete the square

$$\text{to have } Q = -\frac{1}{2} (\theta - A^{-1}b)^T A (\theta - A^{-1}b) + \frac{1}{2} b^T A^{-1} b.$$

Then, we have: $\overset{\downarrow}{e^Q} \propto N(A^{-1}b, A^{-1})$

$$P(y|x, D) = \frac{1}{(2\pi)^{n/2} |A|^{1/2}} \exp(Q) d\theta + e^C$$

Now with some factoring we see:

$$P(y|x, D) = \gamma \int_{\theta} N(\theta^{-1}b, A^{-1}) \cdot e^{\frac{1}{2} b^T A^{-1} b} d\theta$$

+ function not function
of θ of θ , constant

where $\gamma = \frac{(2\pi)^{\frac{n}{2}}}{(2\pi)^{\frac{n}{2}} \delta^n |z_0|}$

then,

$$P(y|x, D) = \gamma e^{\frac{1}{2} b^T A^{-1} b} \int_{\theta} \underbrace{N(\theta^{-1}b, A^{-1}) d\theta}_{\text{gaussian distribution integral} = 1}$$
$$= \gamma e^{\frac{1}{2} b^T A^{-1} b}$$

No where to proceed it should lead to our gaussian function for the posterior predictive distribution

1.

part 3 cont. we see pure integration does not lead to reasonable results - we are stuck.

12

So, we instead use normal distribution properties to compute the posterior predictive distribution:

Our Bayesian regression model:

$$y = X\theta + \epsilon, \epsilon \sim N(0, \frac{\sigma^2}{r} I_r)$$

$$\begin{matrix} y \in \mathbb{R}^n \\ x \in \mathbb{R}^{n \times d} \end{matrix}$$

number of data points.

$$\theta \in \mathbb{R}^d \Rightarrow P(y|x, \theta) = N(\theta^T x, \sigma^2)$$

Our posterior:

$$P(\theta|D) \sim N(\theta_{MAP}, \Sigma_0)$$

$$\text{Since } P(y|x, D) = \int_{\theta} P(y|x, \theta) \cdot P(\theta|D) d\theta$$

or our posterior predictive model is the convolution of 2 gaussians.

So, $P(y|x, D)$ is a gaussian distribution.

Recall the linear transformation rule on a normal distribution?

normal w.r.t y

$$y = Ax, x \sim N(\mu, \Sigma), P(y) = N(y|A\mu, A\Sigma A^T)$$

Proof:

13

$$x \sim N(\mu, \Sigma), \quad E[x] = \mu.$$

$$Y = Ax. \quad E[Y] = E[AX] = A E[X] = A\mu.$$

$$\text{so, mean of } Y, \quad E[Y] = A\mu.$$

$$\text{Cov}(X) = E[(x - E[x])(x - E[x])^T] = \Sigma$$

$$\text{Cov}(Y) = \text{Cov}(AX) = E[(AX - E[AX])(AX - E[AX])^T]$$

$$= E[(AX - A\mu)(Ax - A\mu)^T]$$

$$= E[A(x - \mu) \cdot (A(x - \mu))^T]$$

$$= E[A(x - \mu)(x - \mu)^T A^T]$$

$$= A E[(x - \mu)(x - \mu)^T] A^T$$

$$= A \cdot \Sigma A^T.$$

Now, going back to $P(y^*|x^*, D)$: where x^*, y^* is the new point(s) being predicted:

$$E[y^*|x^*, D] = E[\underbrace{E[y^*|x^*, \theta]}_{\text{from } P(y|x, \theta)} | D]$$

from $P(y|x, \theta)$.

$$= E[(x^*)^T \theta | D]$$

↓

per our model, $P(y|x, \theta)$

$$= (x^*)^T E[\theta | D] = x^{*T} \theta_{MAP}.$$

per prior, $\theta \sim \mathcal{N}(\theta_{MAP}, \Sigma_{\theta})$

Thus, the mean of $P(y^*|x^*, D) = (x^*)^T \theta_{MAP}$.

14

Now, for the posterior predictive variance:

$$\text{Var}(y^*|x^*, D) = \underbrace{\mathbb{E}[\text{Var}(y^*|x^*, \theta)|D]}_{\text{Total Variance}} + \underbrace{\text{Var}(\mathbb{E}[y^*|x^*, \theta]|D)}_{\text{Law}}$$

$$\mathbb{E}[\text{Var}(y^*|x^*, \theta)|D] = \mathbb{E}[\delta^2|D] = \sigma^2$$

$\underbrace{\quad}_{\sigma^2, \text{ per } P(y|x, \theta) \text{ model}}$

$$\text{Var}(\mathbb{E}[y^*|x^*, \theta]|D) = \text{Var}(x^*)^T \theta_{MAP} |D)$$

mean of $P(y^*|x^*, D)$

$$= x^{*T} (\text{cov}(\theta|D) x^*)$$

linear transformation rules as previously explained.

$$= x^{*T} \Sigma_\theta x^*$$

$$\text{Thus, } \text{Var}(y^*|x^*, D) = \sigma^2 + x^{*T} \Sigma_\theta x^*$$

So,

$$P(y^*|x^*, D) = N(x^{*T} \theta_{MAP}, x^{*T} \Sigma_\theta x^* + \sigma^2)$$

or:

$$P(y|x, D) = N(x^T \theta_{MAP}, x^T \Sigma_\theta x + \sigma^2)$$

Law of Total Variance formula found/used from Wikipedia:

https://en.wikipedia.org/wiki/Law_of_total_variance

1.2 - plotting uncertainty band

In [3]:

```
import warnings
warnings.filterwarnings("ignore", category=FutureWarning)
warnings.filterwarnings('ignore')
```

```

import numpy as np
import matplotlib.pyplot as plt
plt.rcParams['figure.figsize'] = [8, 4]
np.random.seed(0)

## please keep the below codes for generating x, DO NOT CHANGE!
# =====
n_train_samples = 6
X_train = np.sort(np.random.rand(n_train_samples))

order_true = 5
true_coefficient = np.array([.5, -1, -.5, -2, 5])

def true_fn(X):
    f = np.ones((X.shape))
    for i in range(order_true):
        f += true_coefficient[i] * X**i
    return f

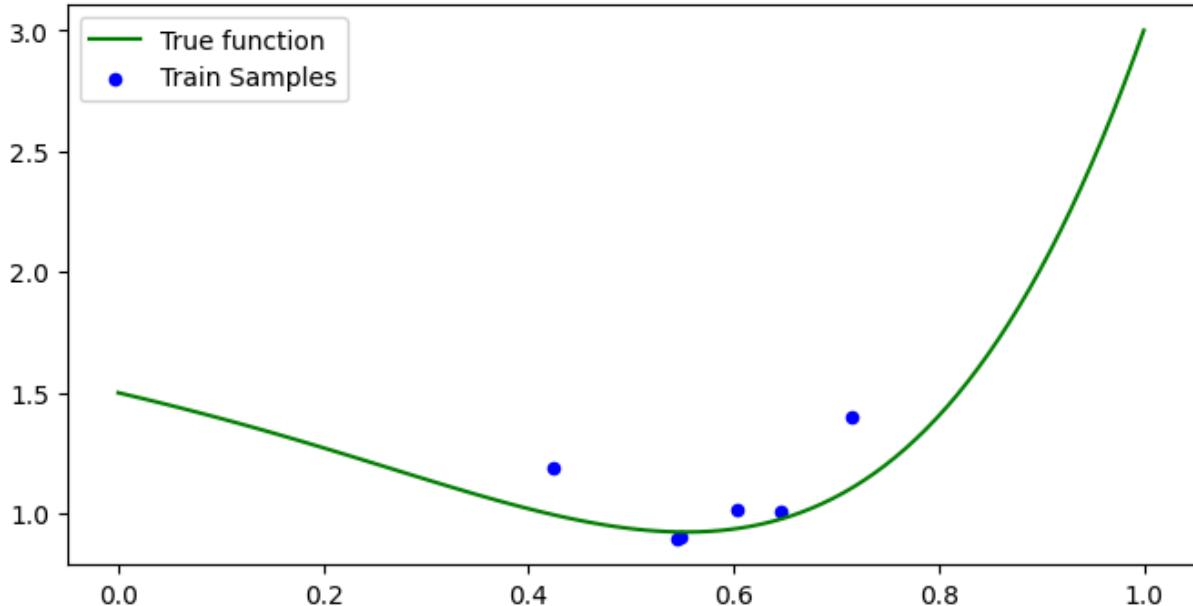
X_test = np.linspace(0., 1., 100)
## adding y_train
y_train = true_fn(X_train) + (np.random.randn(n_train_samples)) * 0.2

plt.plot(X_test,true_fn(X_test),'g',label = 'True function')
plt.scatter(X_train, y_train, edgecolor='b', s=20, facecolor = 'b', label="T
plt.legend()

# [CONTINUE YOUR WORK FROM HERE!]

```

Out[3]: <matplotlib.legend.Legend at 0x7f126c761300>



In [7]:

```

import numpy as np
import matplotlib.pyplot as plt
from sklearn.preprocessing import PolynomialFeatures
from numpy.linalg import inv

```

```

# Add Gaussian noise with delta = 0.2
delta = 0.2
# y_train = true_fn(X_train) + np.random.normal(0, delta, size=X_train.shape)

# Plot the true function and noisy data samples
X_test = np.linspace(0, 1, 100).reshape(-1, 1)
plt.plot(X_test, true_fn(X_test), 'g', label='True function')
plt.scatter(X_train, y_train, edgecolor='b', s=20, facecolor='b', label='Train data')
plt.legend()

# Bayesian Linear Regression with MAP estimate
def design_matrix(X, degree):
    poly = PolynomialFeatures(degree=degree)
    return poly.fit_transform(X.reshape(-1, 1))

degree = 2 # Assuming 2nd order model for PolynomialFeatures
Phi_train = design_matrix(X_train, degree)

# Prior settings (regularization parameter lambda)
lamda = 3 # Prior variance

# Compute theta_MAP (MAP estimate)
I = np.eye(Phi_train.shape[1]) # Identity matrix
theta_map = inv(Phi_train.T @ Phi_train + (delta**2 / lamda**2) * I) @ Phi_train.T @ y_train

# Compute posterior covariance (sigma_0)
sigma_0 = inv((Phi_train.T @ Phi_train)/delta**2 + I/lamda**2)

# Predictive distribution
def predictive_distribution(X_test, theta_map, sigma_0, degree):
    Phi_test = design_matrix(X_test, degree)

    mean_pred = Phi_test @ theta_map
    var_pred = np.sum(Phi_test @ sigma_0 @ Phi_test.T, axis=1) + delta**2
    return mean_pred, np.sqrt(var_pred)

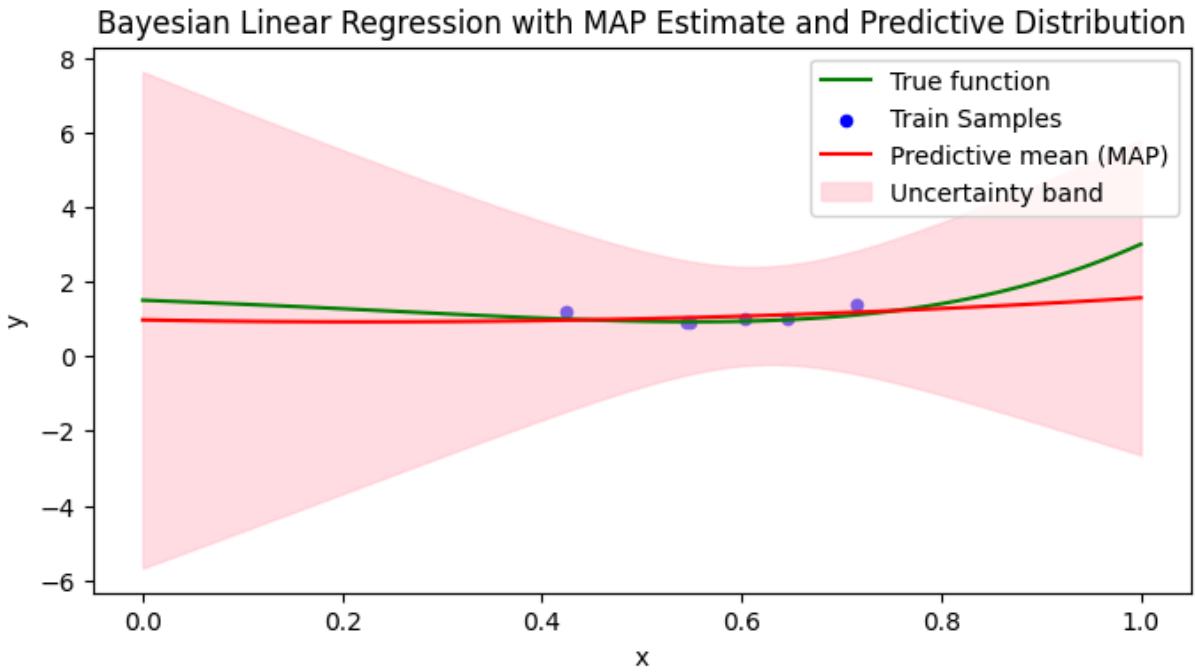
mean_pred, std_pred = predictive_distribution(X_test, theta_map, sigma_0, degree)

# Plot predictive mean and uncertainty bands
plt.plot(X_test, mean_pred, 'r', label='Predictive mean (MAP)')
plt.fill_between(X_test.flatten(), mean_pred.flatten()-2*std_pred.flatten(),
                 mean_pred.flatten()+2*std_pred.flatten(), color='pink', alpha=0.5)

plt.legend()
plt.title('Bayesian Linear Regression with MAP Estimate and Predictive Distribution')
plt.xlabel('x')
plt.ylabel('y')
plt.show()

# Discussion:
# - theta_map is the result of optimizing the regularized loss function.
# - The prior variance ( $\lambda^2$ ) acts as regularization to avoid overfitting.
# - The predictive distribution takes into account the posterior uncertainty

```



1.2 - Notes

The prior variance of 2 ($\lambda^2 = 3$) versus the data's variance of 0.04 ($\delta^2 = 0.04$) is chosen to have a reasonable regularization value. Note that the optimal theta (via the MAP approach) is a least squares solution with regularization = $\delta^2/\lambda^2 = 0.0133$. Also, the uncertainty band is shown with 2 standard deviations.

1.3 - benefits of Bayesian vs Maximum Likelihood

By incorporating prior information of the parameters with the use of prior distributions, the Bayesian methods can take into consideration previous data or domain knowledge to inform current analysis. The maximum likelihood approach does not allow incorporation of prior information or results, and relies solely on the observed data at the current time.

Bayesian methods provide a full posterior distribution (posterior predictive distribution) which provides the complete picture of both the appropriate value and uncertainty for each estimate; meanwhile, maximum likelihood can only provide an estimate of the parameter without quantifying the uncertainty of the estimate.

With the use of the prior, Bayesian methods can be more stable for small sample sizes; with just the small sample size (and no prior), maximum likelihood estimates may produce widely unreliable results.

Question 2 (20 Points) Mean and Covariance

In lecture 7, we maximize the likelihood:

$$\max_{\boldsymbol{\mu}_k, \Sigma_k} \sum_{i:y^{(i)}=k} \log P(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_k, \Sigma_k) = \max_{\boldsymbol{\mu}_k, \Sigma_k} \sum_{i:y^{(i)}=k} \log \mathcal{N}(\mathbf{x}^{(i)} | \boldsymbol{\mu}_k, \Sigma_k).$$

Prove that the empirical means and covariances of each class k are

$$\boldsymbol{\mu}_k = \frac{\sum_{i:y^{(i)}=k} \mathbf{x}^{(i)}}{n_k}$$

$$\Sigma_k = \frac{\sum_{i:y^{(i)}=k} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^\top}{n_k}.$$

where $\mathbf{x}^{(i)} \in \mathbb{R}^n$, $\boldsymbol{\mu} \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$, n is the number of features of training sample $\mathbf{x}^{(i)}$.

2.1 - proof for $\boldsymbol{\mu}_k$

$$2. \text{ Prove: } \mu_k = \frac{\sum_{i:y^i=k} x^i}{n_k}, \quad \Sigma_k = \frac{\sum_{i:y^i=k} (x^i - \mu_k)(x^i - \mu_k)^T}{n_k}$$

for $x^i \in \mathbb{R}^n$, $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$, n is number of features in training sample x^i . (dimension)

Let there be K classes. we maximize:

$$\max_{\mu_k, \Sigma_k} \sum_{i:y^i=k} \log P(x^i | y^i, \mu_k, \Sigma_k) = \max_{\mu_k, \Sigma_k} \sum_{i:y^i=k} \log (N(x^i | \mu_k, \Sigma_k)) = \max_{\mu_k, \Sigma_k} L(\mu_k, \Sigma_k)$$

let n_k be the number of data samples (x^i, y^i) such that $y^i = k$
class $k, k \in [1, K]$

$$N(x^i | \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{n/2} |\Sigma_k|^{1/2}} \exp \left(-\frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k) \right)$$

$$\log N(x^i | \mu_k, \Sigma_k) = -\frac{n}{2} \log (2\pi) - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)$$

Thus, our problem likelihood $L(\mu_k, \Sigma_k)$ is $(x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)$

$$L(\mu_k, \Sigma_k) = -\frac{n_k}{2} \log (2\pi) - \frac{n_k}{2} \log |\Sigma_k| - \sum_{i:y^i=k} \frac{1}{2} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)$$

To maximize $L(\mu_k, \Sigma_k)$, we set $\frac{\partial L}{\partial \mu_k}, \frac{\partial L}{\partial \Sigma_k} = 0$.

$$1) \frac{\partial L}{\partial \mu_k} = \frac{1}{2} \sum_{i:y^i=k} \frac{\partial}{\partial \mu_k} (x^i - \mu_k)^T \Sigma_k^{-1} (x^i - \mu_k)$$

we know $\frac{d}{dx} a^T A x = 2A x$.

16

(shown in HW1).

So, with the chain rule,

$$\frac{\partial L}{\partial \mu_k} = \frac{1}{2} \sum_{i:y^i=k} \Sigma_k^{-1} (x^i - \mu_k) (-2)$$

$$= \sum_{i:y^i=k} \Sigma_k^{-1} (x^i - \mu_k)$$

$$= \Sigma_k^{-1} \sum_{i:y^i=k} (x^i - \mu_k)$$

$$= \Sigma_k^{-1} \left(-n_k \mu_k + \sum_{i:y^i=k} x^i \right)$$

Setting $\frac{\partial L}{\partial \mu_k} = 0$, we have:

$$\Sigma_k^{-1} \left(\sum_{i:y^i=k} x^i - n_k \mu_k \right) = \vec{0}$$

Since Σ_k^{-1} is invertible $(\Sigma_k^{-1})^{-1} = \Sigma_k$ exists &

Σ_k^{-1} is not a null matrix ($\Sigma_k \Sigma_k^{-1} = I$, not null),

we have $\underbrace{\Sigma_k \Sigma_k^{-1}}_I \left(\sum_{i:y^i=k} x^i - n_k \mu_k \right) = \Sigma_k \cdot \vec{0}$

Thus, $\sum_{i:y^i=k} x^i - n_k \mu_k = \vec{0}$.

$$\text{So, } \boxed{\mu_k = \frac{\sum_{i:y^i=k} x^i}{n_k}}$$

2.2 - proof for Σ_k

Now, setting $\frac{\partial L}{\partial \Sigma_k} = 0$:

$$\frac{\partial L}{\partial \Sigma_k} =$$

$$\frac{\partial}{\partial \Sigma_k} \left(\frac{n_k}{2} \log |\Sigma_k| - \frac{1}{2} \sum_{i:y_i=k} (x^i - \mu_k)^\top \Sigma_k^{-1} (x^i - \mu_k) \right)$$

$$\frac{\partial}{\partial \Sigma_k} \frac{n_k}{2} \log |\Sigma_k| = \frac{n_k}{2} \frac{\partial}{\partial \Sigma_k} \log |\Sigma_k|$$

$$\frac{\partial}{\partial \Sigma_k} \log |\Sigma_k| = \frac{1}{|\Sigma_k|} \cdot \frac{\partial}{\partial \Sigma_k} |\Sigma_k|$$

Math Derivations:

18

We know Σ_k is a square, invertible matrix.

We know $C_{ij} = (-1)^{i+j} M_{ij}$; $M_{ij} = \det(\Sigma_{ij})$

where C is the cofactor at row i , column j & M is the minor.

$$\text{Then, } |\Sigma_k| = \sum_{j=1}^n \sum_{ij} C_{ij}$$

$$= \sum_{j=1}^n \sum_{ij} C_{ij} \quad \checkmark \text{ for any } l \in [1, n]$$

let $l = a$

$$\frac{\partial |\Sigma|}{\partial \Sigma_{ab}} = \frac{\partial}{\partial \Sigma_{ab}} \left(\sum_{j=1}^n \sum_{aj} C_{aj} \right)$$

$$= \sum_{j=1}^n \frac{\partial}{\partial \Sigma_{ab}} \sum_{aj} C_{aj}$$

$$= C_{ab} = C_{ba} \rightarrow i+j = a+b \text{ for } C_{ab} \& C_{ba}$$

$M_{ij} = M_{ji}$ since Σ is symmetric.

$$|\Sigma| = \Sigma \cdot C \quad \Sigma \cdot \Sigma^{-1} = I_d$$

$$\Sigma = \frac{1}{C}$$

$$\Sigma^{-1} = \frac{I_d}{\Sigma} = \frac{C}{|\Sigma|}; C = |\Sigma| \Sigma^{-1}$$

$$\text{so, } \frac{\partial |\Sigma|}{\partial \Sigma} = C = |\Sigma| \Sigma^{-1}.$$

$$\text{or, } \frac{\partial}{\partial \Sigma_k} |\Sigma_k| = |\Sigma_k| \Sigma_k^{-1}$$

$$\text{so, } \frac{\partial}{\partial \Sigma_k} \log |\Sigma_k| = \frac{1}{|\Sigma_k|} |\Sigma_k| \Sigma_k^{-1} = \Sigma_k^{-1}$$

19

$$F \text{ Aside: } \frac{\partial}{\partial A} A^{-1} = ?$$

$$A^{-1} A = I \Rightarrow \frac{\partial}{\partial A} (A^{-1} A) = 0.$$

w/ product rule: $\frac{\partial}{\partial A}$

$$\frac{\partial}{\partial A} A^{-1} \cdot A + \frac{\partial}{\partial A} A \cdot A^{-1} = 0$$

$$\frac{\partial}{\partial A} A^{-1} \cdot A + A^{-1} = 0 \Rightarrow \frac{\partial}{\partial A} A^{-1} A = -A^{-1}$$

$$\frac{\partial}{\partial A} A A^{-1} = -A^{-1} A^{-1}; \quad \frac{\partial}{\partial A} A^{-1} = -A^{-1} A^{-1}$$

$$\frac{\partial}{\partial A} v^T A^{-1} v = ? \quad \text{let } f(A) = v^T A^{-1} v \\ \uparrow \text{vector}$$

$$\frac{\partial}{\partial A} f(A) = v^T \frac{\partial}{\partial A} A^{-1} v = v^T (-A^{-1} \cdot A^{-1}) v$$

$$\frac{\partial}{\partial A} (v^T A^{-1} v) = -v^T A^{-1} A^{-1} v$$

$$\frac{\partial L}{\partial \Sigma_k} = -\frac{n_K}{2} \frac{\partial}{\partial \Sigma_k} (\log |\Sigma_k|) - \frac{1}{2} \sum_{i:y_i=K} \frac{\partial}{\partial \Sigma_k} (x_i^T \mu_K)^T \Sigma_k^{-1} (x_i^T \mu_K)$$

$$\frac{\partial L}{\partial \Sigma_k} = -\frac{n_K}{2} \Sigma_k^{-1} + \frac{1}{2} \sum_{i:y_i=K} (x_i^T \mu_K)^T \Sigma_k^{-1} \Sigma_k^{-1} (x_i^T \mu_K)$$

f Aside : $v^T A^{-1} A^{-1} v = A^{-1} v v^T A^{-1}$ when A is symmetric?
 $\hookrightarrow v \in \mathbb{R}^d, A \in \mathbb{R}^{d \times d}$

$$\text{let } p = A^{-1}v, p^T = v^T(A^{-1})^T$$

for $A = A^T$:

$$A^{-1}A = I;$$

$$I = (A^T)^{-1}(A^T) = (A^T)^{-1}A = I; (A^T)^{-1} = A^{-1}$$

$$\text{So, } p^T = v^T A^{-1}$$

$$\text{Thus, } v^T A^{-1} A^{-1} v = p^T p = \sum_{i=1}^d (p_i)^2 = \|A^{-1}v\|^2$$

Now, we use proof by contradiction:

we evaluate the RHS of the hypothesis:

$$A^{-1}v \cdot v^T A^{-1} = p \cdot p^T = \sum_{i=1}^d (p_i)^2 = \|A^{-1}v\|^2$$

Thus for the hypothesis to be false we must have $\|A^{-1}v\|^2 \neq \|A^{-1}v\|^2$ which is obviously false.

Both sides can be simplified to the same value so

$$v^T A^{-1} A^{-1} v = A^{-1} v v^T A^{-1} \quad \square$$

Thus, we have, setting $\frac{\partial L}{\partial \Sigma_k} = 0$,

$$\text{scalar } \frac{n_k}{2} \Sigma_k^{-1} = \frac{1}{2} \sum_{i:y_i=k} (x^i - \mu_k)^T \Sigma_k^{-1} \Sigma_k (x^i - \mu_k)$$

$$\Rightarrow \frac{n_k}{2} \Sigma_k^{-1} = \frac{1}{2} \sum_{i:y_i=k} \Sigma_k^{-1} (x^i - \mu_k)^T (x^i - \mu_k) \Sigma_k$$

$$\Sigma_k^{-1} \frac{n_k}{2} = \frac{1}{2} \sum_{i:y^i=k} \Sigma_k^{-1} (x^i - \mu_k) (x^i - \mu_k)^T \Sigma_k^{-1}$$

multiplying both sides by Σ_k from the left:

$$\Sigma_k \Sigma_k^{-1} \frac{n_k}{2} = \frac{1}{2} \sum_{i:y^i=k} \Sigma_k \Sigma_k^{-1} (x^i - \mu_k) (x^i - \mu_k)^T \Sigma_k^{-1}$$

$$n_k = \sum_{i:y^i=k} (x^i - \mu_k) (x^i - \mu_k)^T \Sigma_k^{-1}$$

Note that since Σ_k is symmetric we can multiply Σ_k from the right on equations.

In the above equation we have a scalar on the left side so, this property isn't necessarily required

$$n_k \Sigma_k = \sum_{i:y^i=k} (x^i - \mu_k) (x^i - \mu_k)^T \Sigma_k^{-1} \Sigma_k$$

\Rightarrow

$$\boxed{\Sigma_k = \frac{\sum_{i:y^i=k} (x^i - \mu_k) (x^i - \mu_k)^T}{n_k}}$$

Question 3: (25 Points) Generative VS Discriminative Models

We are going to use the Iris flower dataset. We will use three features sepal length (cm), sepal width (cm), petal length (cm) to classify types of flowers in generative model.

1. Compute the empirical means and covariance matrices for each type of flowers.
2. Generating new flowers, denoted as set S , from the learned probability and plotting these generated flowers in 3D. Observe and discuss your results.
3. We assume that set S is test data. Classify the test data S using the generative model that you have found in question 4.1. Discuss your results.
4. We assume that set S is test data. Using logistic regression to training the model, then verify the learn model with test data S . Compare and discuss your results with question 4.3

```
In [141...]  

import numpy as np  

import pandas as pd  

import warnings  

warnings.filterwarnings('ignore')  

from sklearn import datasets  
  

# Load the Iris dataset  

iris = datasets.load_iris(as_frame=True)  
  

# print part of the dataset  

iris_X, iris_y = iris.data, iris.target  

pd.concat([iris_X, iris_y], axis=1).head()  
  

#print(iris_X.shape, iris_y.shape)
```

Out[141...]

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

3.0 - data

Note that the different colors for each cluster in the plots is representative of a different iris flower type (i.e., setosa, versicolor, etc.).

```
In [142...]  

# Create a DataFrame for easier handling  

# Select only the desired columns: sepal length, sepal width, and petal length  

iris_X = iris_X[['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)']]  
  

# Rename the columns  

iris_X.columns = ['sepal_length', 'sepal_width', 'petal_length']  

iris_y.name = 'target'  
  

df = pd.concat([iris_X, iris_y], axis=1)  

df.columns = ['sepal_length', 'sepal_width', 'petal_length', 'target']  
  

# Display the first few rows of the dataset  

print(df.head())  

# compare generate vs train  

fig = plt.figure(figsize=(10, 7))  

ax = fig.add_subplot(111, projection='3d')  

ax.set_xlabel('Sepal Length (cm)')  

ax.set_ylabel('Sepal Width (cm)')  

ax.set_zlabel('Petal Length (cm)')  

plt.title('Original Flowers (to be used as testing)')  

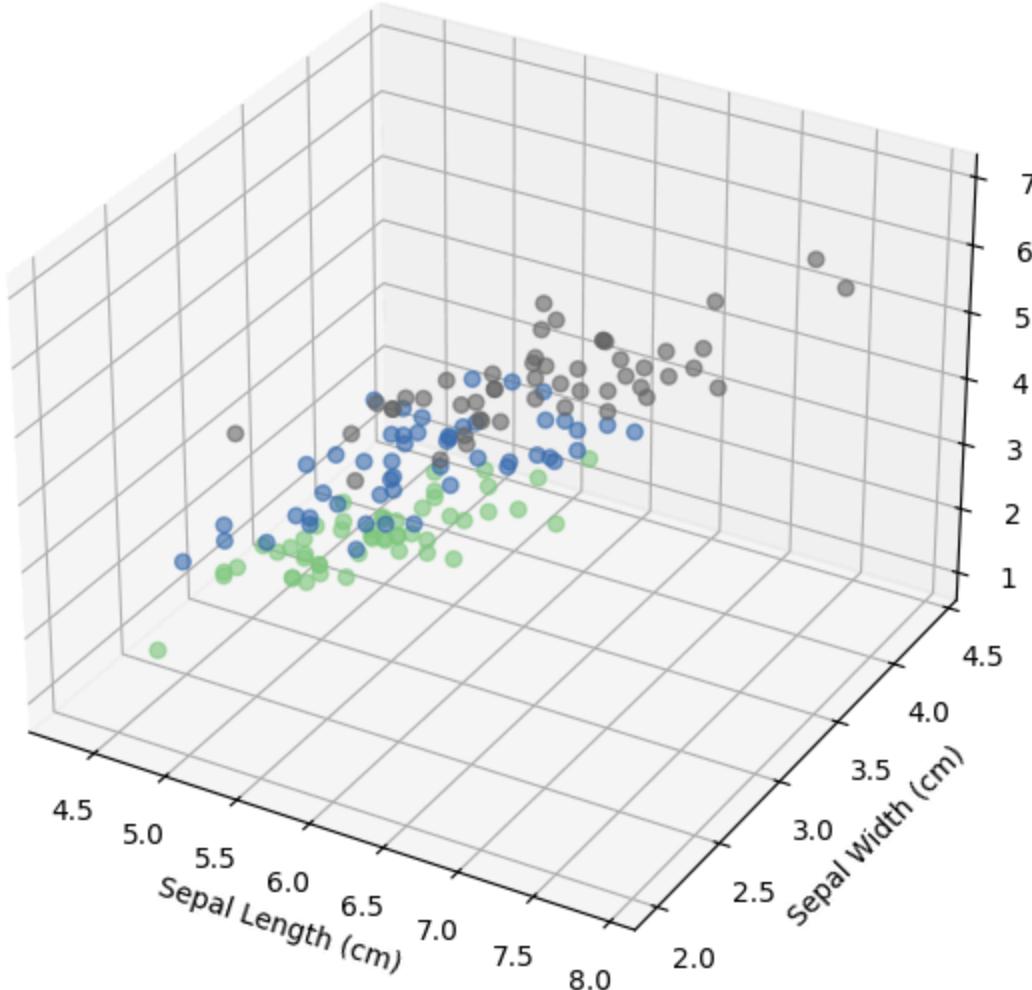
scatter1 = ax.scatter(iris_X['sepal_length'], iris_X['sepal_width'], iris_X[
```

```
ax.legend()  
plt.show()
```

	sepal_length	sepal_width	petal_length	target
0	5.1	3.5	1.4	0
1	4.9	3.0	1.4	0
2	4.7	3.2	1.3	0
3	4.6	3.1	1.5	0
4	5.0	3.6	1.4	0

Original Flowers (to be used as testing)

● Original Flowers



3.1 - mean/covariance for each flower

In [147...]

```
from mpl_toolkits.mplot3d import Axes3D  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score  
  
# 1. Compute the empirical means and covariance matrices for each type of flower  
means = {}
```

```

covariances = {}

for target in np.unique(iris_y):
    flower_data = df[df['target'] == target][['sepal_length', 'sepal_width']]
    means[target] = flower_data.mean().values
    covariances[target] = flower_data.cov().values

print("Empirical Means:")
for target, mean in means.items():
    print(f"Class {target}: {mean}")

print("\nCovariance Matrices:")
for target, cov in covariances.items():
    print(f"Class {target}:\n{cov}\n")

```

Empirical Means:
 Class 0: [5.006 3.428 1.462]
 Class 1: [5.936 2.77 4.26]
 Class 2: [6.588 2.974 5.552]

Covariance Matrices:
 Class 0:
 [[0.12424898 0.09921633 0.0163551]
 [0.09921633 0.1436898 0.01169796]
 [0.0163551 0.01169796 0.03015918]]

 Class 1:
 [[0.26643265 0.08518367 0.18289796]
 [0.08518367 0.09846939 0.08265306]
 [0.18289796 0.08265306 0.22081633]]

 Class 2:
 [[0.40434286 0.09376327 0.3032898]
 [0.09376327 0.10400408 0.07137959]
 [0.3032898 0.07137959 0.30458776]]

3.2 - generate new flowers

```

In [163]: # 2. Generate new flowers from the learned probability and plot them in 3D.
num_samples = 100
S = []

for target in np.unique(iris_y):
    mean = means[target]
    cov = covariances[target]
    # Generate samples from a multivariate normal distribution
    new_samples = np.random.multivariate_normal(mean, cov, num_samples)
    S.append(new_samples)

# Convert the list of arrays to a single array
S = np.vstack(S)
S_labels = np.concatenate([[target]*num_samples for target in np.unique(iris_y)])

```

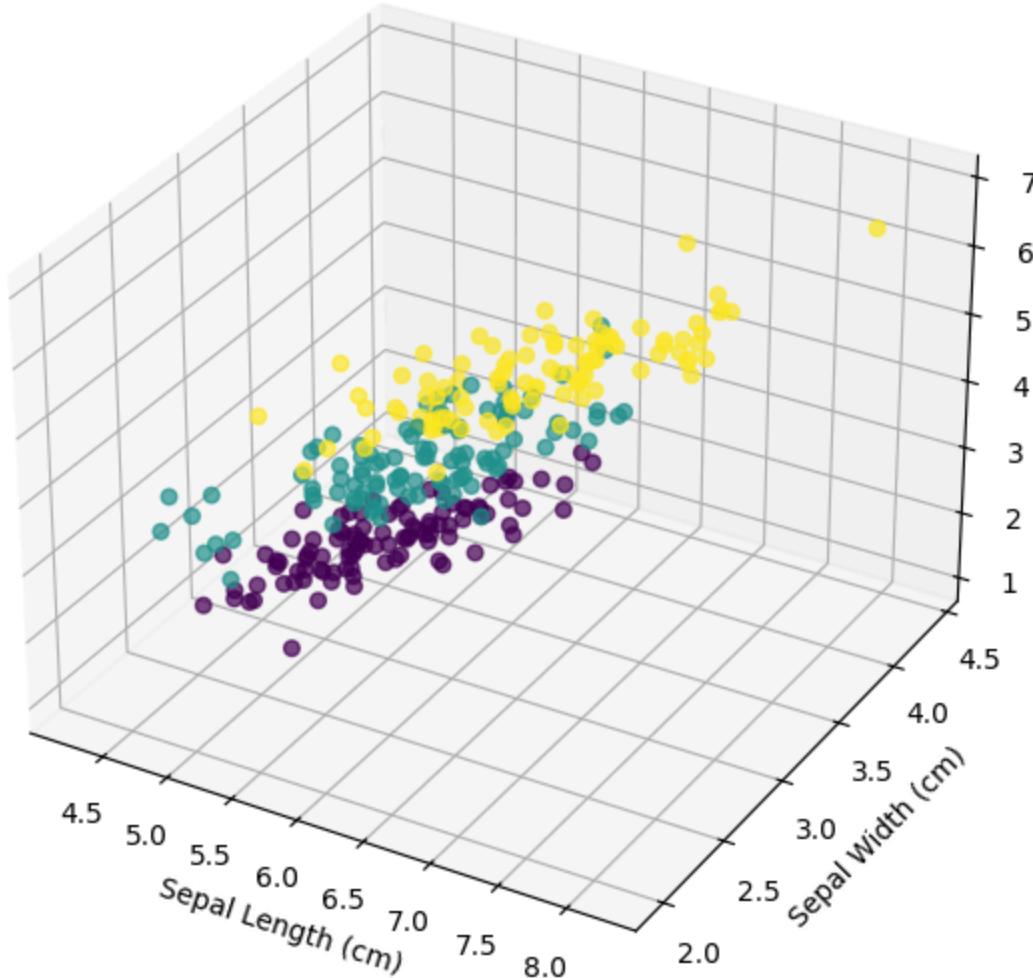
```

fig = plt.figure(figsize=(10, 7))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(S[:, 0], S[:, 1], S[:, 2], c=S_labels, cmap='viridis', s=30, alpha=0.8)
ax.set_xlabel('Sepal Length (cm)')
ax.set_ylabel('Sepal Width (cm)')
ax.set_zlabel('Petal Length (cm)')
plt.title('Generated Flowers from the Generative Model')

```

Out[163... Text(0.5, 0.92, 'Generated Flowers from the Generative Model')

Generated Flowers from the Generative Model



In [167... # compare generate vs train

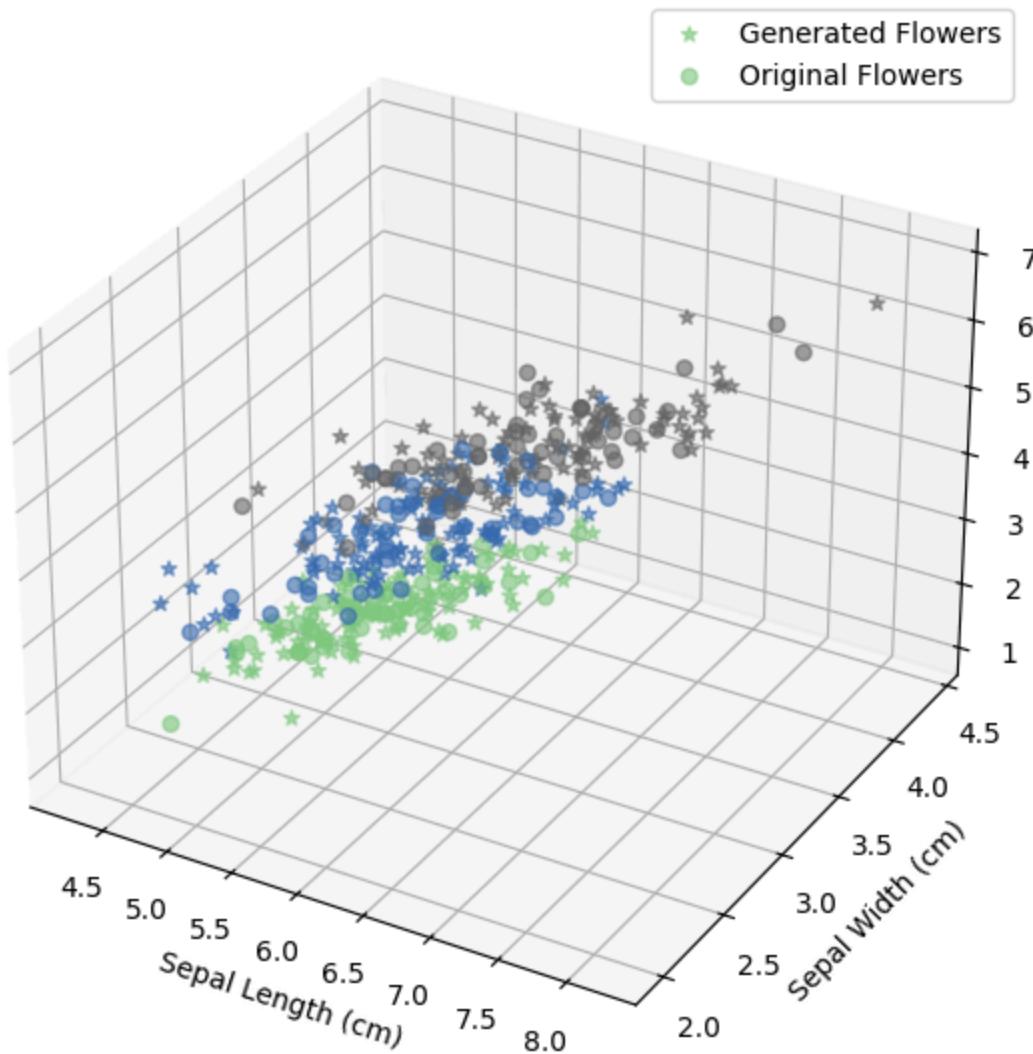
```

fig = plt.figure(figsize=(10, 7))
ax = fig.add_subplot(111, projection='3d')
ax.scatter(S[:, 0], S[:, 1], S[:, 2], c=S_labels, cmap='Accent', s=30, alpha=0.8)
ax.set_xlabel('Sepal Length (cm)')
ax.set_ylabel('Sepal Width (cm)')
ax.set_zlabel('Petal Length (cm)')
plt.title('Generated Flowers from the Generative Model')
scatter1 = ax.scatter(iris_X['sepal_length'], iris_X['sepal_width'], iris_X['petal_length'], c=iris_y, cmap='viridis', s=30, alpha=0.8)

```

```
ax.legend()  
plt.show()
```

Generated Flowers from the Generative Model



Discussion

Generated Flowers: The new flowers were generated using a multivariate normal distribution based on the empirical means and covariance matrices of the original Iris dataset. This process simulates how new data points might look if they follow the same statistical properties as the original data.

3D Plot: The plot shows the generated flowers in a 3D space defined by sepal length, sepal width, and petal length. Each point is colored according to its class label, providing a visual representation of the distribution of the generated flowers.

Distribution: The generated flowers should ideally follow the same distribution as the original flowers. If the generative model has captured the underlying data distribution well, the generated points will be spread out in a similar manner to

the original data points. This is especially obvious in the second plot, where the true data and the generated data are overlayed/Plotted on the same plot. Note that the different colors are representative of different Iris groups; the similar groupings of the generated and true data are obvious, indicating a well-built generative model.

Model Validation: This visualization helps validate the generative model by comparing the spread and clustering of the generated flowers to the original dataset. If the generated flowers appear to be well-distributed and form similar clusters, it indicates that the model has effectively learned the underlying distribution of the data (this is true, per our visual observation).

3.3 - classify test data S

```
In [168]: # 3. Classify the test data using the generative model.  
# We'll use the multivariate Gaussian likelihood for classification  
def classify_generative(samples, means, covariances):  
    predictions = []  
    for sample in samples:  
        probs = []  
        for target in np.unique(iris_y):  
            mean = means[target]  
            cov = covariances[target]  
            # Calculate the multivariate Gaussian probability density  
            diff = sample - mean  
            prob = np.exp(-0.5 * diff.T @ np.linalg.inv(cov) @ diff) / np.sqrt(np.linalg.det(cov))  
            probs.append(prob)  
        predictions.append(np.argmax(probs))  
    return np.array(predictions)  
  
predictions_generative = classify_generative(S, means, covariances)  
accuracy_generative = accuracy_score(S_labels, predictions_generative)  
print(f"Generative Model Classification Accuracy: {accuracy_generative:.4f}")
```

Generative Model Classification Accuracy: 0.9467

Discussion

The generative model accuracy is $\approx 94.67\%$. The high accuracy is indicative of an effective generative model (effective at classifying test data). This suggests that the model has successfully captured the underlying data distribution and can distinguish between different classes.

Note that the generative model assumes the features follow a gaussian distribution; the high accuracy implies that this assumption is reasonably valid for the Iris dataset, which led the model to perform well.

3.4 - logistic regression

```
In [169]: # Fit logistic regression
logistic_model = LogisticRegression(max_iter=500)
logistic_model.fit(iris_X, iris_y)

# Verify the learned model with test data S
predictions_logistic = logistic_model.predict(S)
accuracy_logistic = accuracy_score(S_labels, predictions_logistic)
print(f"Logistic Regression Classification Accuracy: {accuracy_logistic:.4f}")
```

Logistic Regression Classification Accuracy: 0.9400

Discussion

The discriminative logistic regression model accuracy is $\approx 94.00\%$. The high accuracy is indicative of an effective discriminative model (effective at classifying test data). This suggests that the model has successfully captured the underlying data distribution and can distinguish between different classes. Note that discriminative models directly model $P(y|x)$ and do not need assumptions about the gaussian features. So, we would expect the logistic regression to be at a higher accuracy than the generative model. Generally, since discriminative models don't model (or assume anything about) the joint distribution $P(y, x|\theta)$ or $P(x|\theta)$, discriminative models are more accurate for classification than compared with generative models. In this case, we can visually see that the each feature follows an approximate gaussian distribution. Thus, this is a type of prior information that allows GDA to be a better model than the discriminative logistic regression model. Also, since logistic regression only allows linear boundaries, and since the iris dataset is obviously not linearly separable, GDA, which allows for quadratic boundaries, can be a more appropriate classification model.

Notes about different probability distributions:

$P(x|y)$: This is the probability of the features (x) given the class label (y). Generative models model this conditional probability. They assume a specific distribution for the features given the class (e.g., Gaussian distribution).

$P(y)$: This is the prior model probability of the class label (y). Generative models also model/postulate this probability.

$P(x)$: This is the marginal probability of the features (x). It can be computed using the law of total probability: ($P(x) = \sum_y P(x|y)P(y)$), but is generally not needed for analysis (acts as a constant).

$P(y|x)$: This is the posterior probability of the class label (y) given the features (x). Generative models use Bayes' theorem to compute this. This is the main purpose of a generative model.

Question 4: (25 Points) Bernoulli Naive Bayes Model

1. Derive the formula, from lecture 8, for the optimal value of ψ_{jk}

$$\psi_{jk} = \frac{n_{jk}}{n_k}.$$

2. In lecture 8, we classified a document by the Naive Bayes (bag of words) method. You will now make a simple comparison between Naive Bayes and Logistic Regression. Remove class 2 and 3 from the training data. By doing so, our data becomes a binary classification problem. Perform the Naive Bayes (bag of words) method and logistic regression method for the obtained 2-class data. Discuss the results.

4.1 - optimal ψ_{jk}

i) Derive $\psi_{jk} = \frac{n_{jk}}{n_k}$

$\psi_{jk} = P(x=1 | y=k)$, probability of seeing word that belongs to vocabulary set $\phi(x)$, given restriction to category k (e.g. spam email category).

n_{jk} : frequency of word j in class k

n_k : number of documents in class k .

for dataset $D = \{(x^i, y^i) | i=1, 2, \dots, n\}$ we optimize/maximize

the log-likelihood $\log P_\theta(x, y)$ w/ maximum likelihood:

$$\theta = \{\phi_1, \phi_2, \dots, \phi_K, \psi_{11}, \psi_{12}, \dots, \psi_{dK}\}$$

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P_\theta(x^i, y^i) = \underset{\theta}{\operatorname{argmax}} l(\theta)$$

$$P_\theta(x, y) = P_\theta(x|y) \cdot P_\theta(y).$$

Note:

$$P_\theta(x; = 1 | y=k) = \text{Bernoulli}(\psi_{jk})$$

so,

$$P_\theta(y) = \text{categorical}(\phi_1, \dots, \phi_K)$$

$$\begin{aligned} l(\theta) &= \sum_{i=1}^n \underbrace{\log P_\theta(x^i | y^i)}_{\substack{\text{only dependent} \\ \text{on } \vec{\psi}}} + \sum_{i=1}^n \underbrace{\log P_\theta(y^i)}_{\substack{\text{only dependent on } \vec{\phi}}} \\ &= \sum_{k=1}^K \sum_{j=1}^d \sum_{i: y^i=k} \log P(x^{(i)} | \psi_{jk}) + \sum_{i=1}^n \log P(y^{(i)} | \vec{\phi}) \end{aligned}$$

we know:

$$P(x_j^{(i)} | \psi_{jk}) = \psi_{jk}^{x_j^{(i)}} (1 - \psi_{jk})^{1-x_j^{(i)}}$$

so,

$$L(\theta) = \sum_{k=1}^K \sum_{j=1}^d \sum_{i:y^i=k} \left(x_j^{(i)} \log \psi_{jk} + (1-x_j^{(i)}) (1 - \psi_{jk}) \right)$$

$$+ \sum_{i=1}^n \log P(y^i | \vec{\phi})$$

to find the maxima of $L(\theta)$, we set $\frac{\partial L}{\partial \theta} = 0$ &
to get $\frac{\partial L}{\partial \psi_{jk}}$

$$\frac{\partial L}{\partial \theta} = 0.$$

$$1) \frac{\partial L}{\partial \psi_{jk}} = \sum_{i:y^i=k} \frac{\partial}{\partial \psi_{jk}} \left[x_j^{(i)} \log \psi_{jk} + (1-x_j^{(i)}) (1 - \psi_{jk}) \right] = 0$$

$$\frac{\partial L}{\partial \psi_{jk}} = \sum_{i:y^i=k} \frac{1}{\psi_{jk}} x_j^{(i)} - \frac{1}{1-\psi_{jk}} (1-x_j^{(i)})$$

$$\Rightarrow \sum_{i:y^i=k} \frac{x_j^{(i)}}{\psi_{jk}} - \frac{1-x_j^{(i)}}{1-\psi_{jk}} = 0.$$

$$\sum_{i:y^i=k} \frac{x_j^{(i)}}{\psi_{jk}} = \sum_{i:y^i=k} \frac{1-x_j^{(i)}}{1-\psi_{jk}}$$

using our definition for n_{jk} and rearranging:

$$\sum_{i:y^{(i)}=k} x_j^{(i)} (1 - \psi_{jk}) = \sum_{i:y^{(i)}=k} (1 - x_j^{(i)}) \psi_{jk}$$

24

recall that n_{jk} is the number of times $x_j^{(i)} = 1$
(vs. $x_j^{(i)} = 0$).

$n_k = \sum n_{jk}$, total # of observations of class k.

so, we have:

$$n_{jk} (1 - \psi_{jk}) = (n_k - n_{jk}) \psi_{jk}$$

$$\Rightarrow n_{jk} = n_k \psi_{jk} \quad \text{or} \quad \psi_{jk} = \frac{n_{jk}}{n_k}$$

4.2 - Naive Bayes

In [89]:

```
import numpy as np
# import pandas as pd
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import CountVectorizer

# for this lecture, we will restrict our attention to just 4 different newsgroups
categories = ['alt.atheism', 'soc.religion.christian', 'comp.graphics', 'sci.space']

# load the dataset
twenty_train = fetch_20newsgroups(subset='train', categories=categories, shuffle=True)

# vectorize the training set
count_vect = CountVectorizer(binary=True, max_features=1000)
y_train = twenty_train.target
X_train = count_vect.fit_transform(twenty_train.data).toarray()
X_train.shape

print('NOTE: our data has 4 class group')
print(y_train)
print('You NEED to remove class 2 and 3 to create the 2-class data')
y_train
```

NOTE: our data has 4 class group

[1 1 3 ... 2 2 2]

You NEED to remove class 2 and 3 to create the 2-class data

Out[89]: array([1, 1, 3, ..., 2, 2, 2])

In [92]:

```
import numpy as np
# import pandas as pd
from sklearn.datasets import fetch_20newsgroups
from sklearn.feature_extraction.text import CountVectorizer
```

```

from sklearn.naive_bayes import MultinomialNB
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, classification_report

twenty_test = fetch_20newsgroups(subset='test', categories=categories, shuffle=True)

# vectorize the training set
count_vect_test = CountVectorizer(binary=True, max_features=1000)
y_test = twenty_test.target
X_test = count_vect_test.fit_transform(twenty_test.data).toarray()

# Filter out classes 2 and 3
mask_train = np.isin(y_train, [0, 1]) # Keep only classes 0 and 1
X_train_split = X_train[mask_train]
y_train_split = y_train[mask_train]
mask_test = np.isin(y_test, [0, 1]) # Keep only classes 0 and 1
X_test_split = X_test[mask_test]
y_test_split = y_test[mask_test]
# # Split the data into training and testing sets
# X_train_split, X_test_split, y_train_split, y_test_split = train_test_split(
#     X_train_binary, y_train_binary, test_size=0.2, random_state=42
# )

# Train Naive Bayes
nb_classifier = MultinomialNB()
nb_classifier.fit(X_train_split, y_train_split)
y_pred_nb = nb_classifier.predict(X_test_split)
y_pred_train_nb = nb_classifier.predict(X_train_split)
# Train Logistic Regression
logistic_classifier = LogisticRegression()
logistic_classifier.fit(X_train_split, y_train_split)
y_pred_logistic = logistic_classifier.predict(X_test_split)
y_pred_train_logistic = logistic_classifier.predict(X_train_split)

print("TRAINING RESULTS")

print("Naive Bayes Results:")
print("Accuracy:", accuracy_score(y_train_split, y_pred_train_nb))
print(classification_report(y_test_split, y_pred_nb, target_names=twenty_train_labels))

print("Logistic Regression Results:")
print("Accuracy:", accuracy_score(y_train_split, y_pred_train_logistic))
print(classification_report(y_test_split, y_pred_logistic, target_names=twenty_train_labels))

print("\nTESTING RESULTS")
# Evaluate and compare results
print("Naive Bayes Results:")
print("Accuracy:", accuracy_score(y_test_split, y_pred_nb))
print(classification_report(y_test_split, y_pred_nb, target_names=twenty_test_labels))

print("Logistic Regression Results:")

```

```

print("Accuracy:", accuracy_score(y_test_split, y_pred_logistic))
print(classification_report(y_test_split, y_pred_logistic, target_names=twer)

```

TRAINING RESULTS

Naive Bayes Results:

Accuracy: 0.9821428571428571

	precision	recall	f1-score	support
alt.atheism	0.52	1.00	0.69	319
comp.graphics	1.00	0.26	0.41	389
accuracy			0.59	708
macro avg	0.76	0.63	0.55	708
weighted avg	0.79	0.59	0.53	708

Logistic Regression Results:

Accuracy: 1.0

	precision	recall	f1-score	support
alt.atheism	0.81	0.77	0.79	319
comp.graphics	0.82	0.85	0.84	389
accuracy			0.82	708
macro avg	0.82	0.81	0.82	708
weighted avg	0.82	0.82	0.82	708

TESTING RESULTS

Naive Bayes Results:

Accuracy: 0.5918079096045198

	precision	recall	f1-score	support
alt.atheism	0.52	1.00	0.69	319
comp.graphics	1.00	0.26	0.41	389
accuracy			0.59	708
macro avg	0.76	0.63	0.55	708
weighted avg	0.79	0.59	0.53	708

Logistic Regression Results:

Accuracy: 0.8177966101694916

	precision	recall	f1-score	support
alt.atheism	0.81	0.77	0.79	319
comp.graphics	0.82	0.85	0.84	389
accuracy			0.82	708
macro avg	0.82	0.81	0.82	708
weighted avg	0.82	0.82	0.82	708

4.2 - Discussion

Generally, since discriminative models don't model (or assume anything about) the joint distribution $P(y,x|\theta)$ or $P(x|\theta)$, discriminative models are more

accurate for classification than compared with generative models. Since we have no information (prior) in regards to the newspaper/article classification, the discriminative model will obviously outperform the Naive Bayes model which must assume a gaussian distribution ($P(x|y=Y) \sim \text{Normal}$) and that each feature can independently predict the class.

Comparing the testing results with the training results, as expected, our training results are more accurate; the model is trained to fit the training data, not the testing data, so, obviously, the model is more well-fit to the training than testing data, and thus shows a higher accuracy on the training data.

The 100% accuracy on the logistic regression training data suggests that both classes can be separated (with respect to the training data). The high Naive Bayes training accuracy (98%) suggests that the training data can also be well fit with a quadratic decision boundary. This is obviously true, since, if a linear model is sufficient, any higher dimensional model is unnecessary, yet can still relatively well fit to the training data. The higher dimensional model, in this case, is less well suited to the training data, due to the above reasoning (assumptions about gaussian data distribution).

This notebook was converted with convert.ploomber.io