

Anushi R. Sadam

Professor Tan Bui Thanh

COE 379L: Introduction to Machine Learning & Data Science
29 October 2024

Homework 4

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

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

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

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.

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

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

mean of $P(y|x_i, \theta) = \theta^T x_i$:

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

proven in HW3.

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 $\bar{\theta}$,
variance of each element λ^2 :

(2b)

Now, we assume a prior of mean 0, variance λ^2 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 ($\mu_i = 0, \mu_j = 0$),

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

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

\rightarrow 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,

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

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

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

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

$P(\theta | D)$

we seek to maximize¹ feature
matrix

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

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

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

$$\text{let } A = \left[\frac{x^T x}{s^2} + \frac{I}{\lambda^2} \right] \quad \& b = \frac{x^T y}{s^2} \quad \& c = \frac{y^T y}{s^2}$$

then:

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

F Aside:

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

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

since:

LHS expansion:

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

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

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

$$A^T = ((mX^T X) + nI)^T = (nI)^T + (mX^T X)^T$$

$$= nI^T + m(X^T X)^T = nI + mX^T (X^T)^T \\ = nI + mX^T X = mX^T X + nI = A.$$

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:

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

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

continuing w/ our exponential factoring:

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

then in the exponent, we have

b

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

so:

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

multivariate Gaussian expression

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

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

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

$$\Sigma_0 = \left[\frac{X^T X}{s^2} + \frac{I}{\lambda^2} \right]^{-1} = s^2 \left[X^T X + \frac{s^2}{\lambda^2} I \right]^{-1}$$

① Note about scalar inverse pg 7

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

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

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

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

The solution to the L_2 regularized solution is

loss:

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

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

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

$$\theta^*_{MAP} = (X^T X + \lambda I)^{-1} X^T y.$$

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

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

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 .

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

$$\textcircled{1} \quad (CA)^{-1} = C^{-1} A^{-1} \quad \text{if } C \in \mathbb{R}, A \in \mathbb{R}^{d \times d}, \lambda \in \mathbb{R}$$

$\downarrow A \text{ must be invertible}$

Proof:

Suppose CA has inverse B :

$$(CA)B = I; \quad C(AB) = I; \quad AB = \frac{1}{C} I$$

$$A^{-1}AB = A^{-1}I; \quad IB = \frac{1}{C} A^{-1}$$

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

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

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

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

where d is the dimension of the data

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

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).

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

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.

9

Q3. Now, to find $P(y|x, D)$:

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

$$= \int_{\theta} \prod_{i=1}^n N(\theta^T x_i, \delta) \cdot N(\theta_{MAP}, \Sigma_0) d\theta$$

Note that

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

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

where X is the feature Matrix.

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

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

let M be the exponent term.

m

We expand M to get:

$$-\frac{1}{2\delta^2} (y^T y - 2y^T X \theta + \theta^T X^T X \theta) - \frac{1}{2} (\theta^T \Sigma_0^{-1} \theta - 2\theta^T \Sigma_0^{-1} \theta_{MAP} + \theta_{MAP}^T \Sigma_0^{-1} \theta_{MAP})$$

Grouping quadratic terms:

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

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

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

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

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

be the exponent in the integral,

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

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: $e^Q \propto N(A^{-1}b, A^{-1})$

$$P(y|x, D) = \int \frac{1}{(2\pi)^n s^n |\Sigma_0|} \exp(Q) d\theta \cdot e^C$$

now with some factoring we see:

11

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

γ function of θ not function of θ , constant

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

then,

$$P(y|x, D) = \gamma e^{\frac{1}{2} b^T A^{-1} b} \int_0^\infty N(A^{-1}b, A^{-1}) d\theta$$

gaussian distribution integral = 1.

$$= \gamma e^{\frac{1}{2} b^T A^{-1} b}$$

No where to proceed & 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:

$$y = X\theta + \epsilon, \epsilon \sim N(0, \sigma^2 I_n)$$

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

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

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

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

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

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

$$\begin{aligned} \text{cov}(Y) &= \text{cov}(AX) = E[(AX - E[AX])(AX - E[AX])^T] \\ &= E[(AX - A\mu)(AX - A\mu)^T] \end{aligned}$$

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

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

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

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

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

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

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

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

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

per prior, $\theta \sim \theta_{MAP}$

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

14

Now, for the posterior predictive variance:

$$\text{Var}(y^*|x^*, D) = \underbrace{\mathbb{E}[\text{Var}(y^*|x^*, \theta)|D]}_{\delta^2, \text{ per } P(y|x, \theta) \text{ model}} + \underbrace{\text{Var}(\mathbb{E}[y^*|x^*, \theta]|D)}_{\text{Law}}$$

Total
variance
Law

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

δ^2 , per $P(y|x, \theta)$ model

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

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

$$= x^{*T} \text{Cov}(\theta|D) x^* \quad \begin{matrix} \nearrow & \searrow \\ \text{linear} & \text{transformation} \\ \text{rules as previously} & \text{explained.} \end{matrix}$$

$$= x^{*T} \Sigma_0 x^*$$

$$\text{Thus, } \text{Var}(y^*|x^*, D) = \delta^2 + x^{*T} \Sigma_0 x^*$$

So,

$$P(y^*|x^*, D) = N(x^{*T} \theta_{MAP}, x^{*T} \Sigma_0 x^* + \delta^2)$$

or:

$$P(y|x, D) = N(x^T \theta_{MAP}, x^T \Sigma_0 x + \delta^2)$$

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

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 $(x^i - \mu_k)$

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

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

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

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

(shown in HW1).

So, with the chain rule,

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

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

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

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

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

$$\sum_k^{-1} \left(\sum_{i:y^i=k} \alpha^i - n_k \mu_k \right) \stackrel{!}{=} 0$$

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),

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

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

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

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

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

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

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

Math Derivations:

18

We know Σ_k is a square, invertible matrix.

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

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

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

$$= \sum_{j=1}^n \Sigma_{1j} \frac{C_{1j}}{l_{1j}} + \text{for any } l \in [1, d]$$

let $l = a$

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

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

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

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

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

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

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

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

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

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

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: $\frac{\partial}{\partial A}$

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

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

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

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

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

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

$$\frac{\partial L}{\partial \Sigma_k} = -\frac{n_k}{2} \frac{\partial}{\partial \Sigma_k} (\log |\Sigma_k|) - \frac{1}{2} \sum_{i:y_i=k} \frac{\partial}{\partial \Sigma_k} (x_i^T - \mu_k)^T \Sigma_k^{-1} (x_i^T - \mu_k)$$

$$\frac{\partial L}{\partial \Sigma_k} = -\frac{n_k}{2} \Sigma_k^{-1} + \frac{1}{2} \sum_{i:y_i=k} (x_i^T - \mu_k)^T \Sigma_k^{-1} \Sigma_k^{-1} (x_i^T - \mu_k)$$

20

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

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

for $A = A^T$:

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

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

$$\text{So, } P^T = U^T A^{-1}$$

$$\text{Thus, } U^T A^{-1} A^{-1} U = P^T P = \sum_{i=1}^d (P_i)^2 = \|A^{-1} U\|^2$$

Now, we use proof by contradiction:

we evaluate the RHS of the hypothesis.

$$A^{-1} U \cdot U^T A^{-1} = P \cdot P^T = \sum_{i=1}^d (P_i)^2 = \|A^{-1} U\|^2$$

Thus for the hypothesis to be false we must

have $\|A^{-1} U\|^2 \neq \|A^{-1} V\|^2$ which is obviously false.

Both sides can be simplified to the same value so

$$U^T A^{-1} A^{-1} V = A^{-1} V U^T A^{-1}$$

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

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

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

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

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

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

$$1) \text{ Derive } \psi_{jk} = \frac{n_{jk}}{n_k}$$

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

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

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

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

so,

Note:

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

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

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

We know:

23

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

so,

$$\begin{aligned} L(\phi) &= \sum_{k=1}^K \sum_{j=1}^d \sum_{i:y^i=k} \left(x_j^{(i)} \log \psi_{jk} + (1-x_j^{(i)}) \log (1 - \psi_{jk}) \right) \\ &\quad + \sum_{i=1}^n \log P(y^i | \phi) \end{aligned}$$

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

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

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

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

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

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

using our definition for n_{jk} and rearranging:

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

24

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

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

So, we have:

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

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