

Week I

1. Notable Probability Distributions

.) MULTINOMIAL

Let's start from the simple case of a Bernoulli trial: you have 2 possible outcomes for an experiment / event, which can happen with probability π and $1-\pi$, respectively.

You can think at the same experiment and record it as $(0, 1) / (1, 0)$, i.e. the outcome is one of the two possible.

The multinomial distribution extends this to the case of multiple possible outcomes, i.e. the recorded outcome is one out of the K possible. For example, $K=4$

$(0, 1, 0, 0), (1, 0, 0, 0), \dots, (0, 0, 1, 0) \dots$

The binomial distribution counts the number of "successes" (1) in N independent Bernoulli trials.

The multinomial counts the number y_i in each category i in N independent trials.

Let π_k be the probability of each trial resulting in category k . Let \underline{Y} be the vector of the number of successes in each category over N trials. Then,

$$\underline{Y} \sim \text{Multinomial}(N; \pi_1, \dots, \pi_K) \quad \text{and}$$

the probability density function is

Pdf of Multinomial dist \rightarrow $\left\{ P(\underline{Y}|N, \pi) = \frac{N!}{y_1! y_2! \dots y_K!} \pi_1^{y_1} \pi_2^{y_2} \dots \pi_K^{y_K} \right\}$

Note that the marginal $Y_k \sim \text{Binomial}(N, \pi_k)$

and $\mathbb{E}[Y] = N\pi$

$$\text{Cov}(Y_i, Y_j) = -N\pi_i\pi_j, i \neq j$$

Multivariate Normal (Gaussian)

This is a continuous probability distribution.

$p = \text{dimension of } \underline{Y}$

$$\underline{Y} \sim N_p(\underline{\mu}, \Sigma)$$

mean

\rightarrow covariance: For any linear combination matrix

(non-random) \underline{c} ,

$$\begin{aligned}\text{Var}(\underline{c}^T \underline{Y}) &= \text{Var}(c_1 Y_1 + \dots + c_p Y_p) = \\ &= \underline{c}^T \Sigma \underline{c}\end{aligned}$$

The probability density function is

$$P(\underline{Y} | \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (\underline{Y} - \underline{\mu})^T \Sigma^{-1} (\underline{Y} - \underline{\mu})}$$

Remark: Σ^{-1} is often called Precision matrix.

) The marginal distributions are

$$Y_j \sim N(\mu_j, \sigma_j^2), \text{ where } \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots \\ \sigma_{21} & \dots & \dots \\ \vdots & \dots & \dots \\ \sigma_p^2 & \dots & \dots \end{bmatrix}$$

and $\text{Cov}(Y_i, Y_j) = \sigma_{ij}$

) Any linear combination has normal distribution:

$$\underline{c}^T \underline{Y} \sim N(\underline{c}^T \underline{\mu}, \underline{c}^T \Sigma \underline{c})$$

→ Partitioned Gaussian

Let's now consider a partition $\underline{Y} = \begin{bmatrix} \underline{Y}_1 \\ \underline{Y}_2 \end{bmatrix} \left\{ \begin{array}{l} P_1 \\ P_2 \end{array} \right\}$ dimensions

$$\underline{\mu} = \begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad (\Sigma_{22} = \Sigma_{21}^T)$$

$$\text{and } \Sigma^{-1} = \Lambda = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix}.$$

The marginal density of \underline{Y}_2 :

$$P(\underline{Y}_2) = \int P(\underline{Y}_1, \underline{Y}_2) d\underline{Y}_1$$

$$\text{Let } \underline{\delta}_2 = \underline{Y}_2 - \underline{\mu}_2 = \begin{bmatrix} \delta_{21} \\ \delta_{22} \end{bmatrix}$$

Completing Sq. terms
on this next PG.
Using

$$\text{Now, } -\frac{1}{2} \underline{\delta}_2^T \Lambda \underline{\delta}_2 = -\frac{1}{2} [\delta_{21}^T \delta_{21}] \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix} \begin{bmatrix} \delta_{21} \\ \delta_{22} \end{bmatrix}$$

$$= -\frac{1}{2} \delta_{21}^T \Lambda_{11} \delta_{21} - \frac{1}{2} \delta_{22}^T \Lambda_{22} \delta_{22} - \delta_{21}^T \Lambda_{12} \delta_{22} =$$

$$\textcircled{*} = -\frac{1}{2} \delta_{21}^T \Lambda_{11} \delta_{21} - \frac{1}{2} (\delta_{21} + \Lambda_{11}^{-1} \Lambda_{12} \delta_{21})^T \Lambda_{11} (\delta_{21} + \Lambda_{11}^{-1} \Lambda_{12} \delta_{21})$$

$$+ \frac{1}{2} \delta_{22}^T \Lambda_{22} \delta_{22} + \frac{1}{2} \delta_{21}^T \Lambda_{12} \Lambda_{22}^{-1} \Lambda_{22} \delta_{21}$$

I want to write the second part as a quadratic form in δ_{21}

plus a term that depends only on δ_{22}

(because I am going to integrate w.r.t. δ_{22})

(*) Aside: Completing the square

$$\text{For scalars: } \frac{a}{2}x^2 + bx = \frac{a}{2} \left(x^2 + \frac{2b}{a}x \right) = \\ = \frac{a}{2} \left[\left(x + \frac{b}{a} \right)^2 - \frac{b^2}{a^2} \right]$$

$$\text{For vectors: } (\underline{x} \text{ symmetric}) \quad \frac{1}{2} \underline{x}^T A \underline{x} + \underline{b}^T \underline{x} = \frac{1}{2} (\underline{x} + \underline{m})^T A (\underline{x} + \underline{m}) - c =$$

$$= \frac{1}{2} \underline{x}^T A \underline{x} + \frac{1}{2} \underline{x}^T A \underline{m} + \frac{1}{2} \underline{m}^T A \underline{x} + \frac{1}{2} \underline{m}^T A \underline{m} - c = \\ = \frac{1}{2} \underline{x}^T A \underline{x} + \underline{m}^T A \underline{x} + \underbrace{\frac{1}{2} \underline{m}^T A \underline{m}}_{\text{Set } = 0} - c$$

$$\Rightarrow \underline{b}^T = \underline{m}^T A \Rightarrow \underline{m} = \underline{A}^{-1} \underline{b}$$

$$c = \frac{1}{2} \underline{m}^T A \underline{m} = \frac{1}{2} \underline{b}^T \underline{A}^{-1} \underline{b}$$



$$\Rightarrow \boxed{\frac{1}{2} \underline{x}^T A \underline{x} + \underline{b}^T \underline{x} = \frac{1}{2} \left(\underline{x} + \underline{A}^{-1} \underline{b} \right)^T A \left(\underline{x} + \underline{A}^{-1} \underline{b} \right) + \\ - \frac{1}{2} \underline{b}^T \underline{A}^{-1} \underline{b}}$$

\uparrow Completing Square formula
for vectors

Now let's integrate:

ignore first term of PDF of Normal
& focus on exponent term only.

$$\int p(\underline{y}_2, \underline{y}_1) d\underline{y}_2 \propto \int \exp\left(-\frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22}^{-1} \underline{\delta}_{\underline{y}_2}\right) d\underline{y}_2 =$$

$$= \exp\left(-\frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22} \underline{\delta}_{\underline{y}_2} + \frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22} \Lambda_{22}^{-1} \Lambda_{22} \underline{\delta}_{\underline{y}_2}\right).$$

✓

Term depending
on $\underline{\delta}_{\underline{y}_2}$ comes
out of integral

$$\int \exp\left[-\frac{1}{2} (\underline{\delta}_{\underline{y}_2} + \Lambda_{22}^{-1} \Lambda_{21} \underline{\delta}_{\underline{y}_1})^\top \Lambda_{22} (\underline{\delta}_{\underline{y}_2} + \Lambda_{22}^{-1} \Lambda_{21} \underline{\delta}_{\underline{y}_1})\right] d\underline{y}_2$$

This integral gives
1 but constant comes
out as

$$\propto \overset{\text{darts}}{N} \left(\underline{\mu}_2 - \Lambda_{22}^{-1} \Lambda_{21} \underline{\delta}_{\underline{y}_1}, \Lambda_{22}^{-1} \right)$$

$$\Rightarrow \int p(\underline{y}_2, \underline{y}_1) d\underline{y}_2 \propto (2\pi)^{\frac{1}{2}} |\Lambda_{22}^{-1}|^{\frac{1}{2}} \exp\left(-\frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22} \underline{\delta}_{\underline{y}_2} + \cancel{\frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22} \Lambda_{22}^{-1} \Lambda_{21} \underline{\delta}_{\underline{y}_1}}\right)$$

$$+ \frac{1}{2} \underline{\delta}_{\underline{y}_2}^\top \Lambda_{22} \Lambda_{22}^{-1} \Lambda_{21} \underline{\delta}_{\underline{y}_1})$$

$$\propto \exp\left(-\frac{1}{2} (\underline{y}_2 - \underline{\mu}_2)^\top (\Lambda_{22} - \Lambda_{21} \Lambda_{22}^{-1} \Lambda_{21}^\top) (\underline{y}_2 - \underline{\mu}_2)\right)$$

$$\Rightarrow \underline{y}_2 \sim N\left(\underline{\mu}_2, (\Lambda_{22} - \Lambda_{21} \Lambda_{22}^{-1} \Lambda_{21}^\top)^{-1}\right)$$

\downarrow
Marginal density of \underline{y}_1

→ Written also in next
page using Block inverses.

Ass 2 : Block inverse

Block matrix $\begin{bmatrix} A & B \\ C & D \end{bmatrix}$, A, D square

$$\begin{bmatrix} M & N \\ P & Q \end{bmatrix} = \begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = ?$$

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix} \begin{bmatrix} M & N \\ P & Q \end{bmatrix} = \begin{bmatrix} I & 0 \\ 0 & I \end{bmatrix}$$

$$\Rightarrow AM + BP = I, CM + DP = 0, \Rightarrow DP = -CM$$
$$P = -D^{-1}CM$$

$$AM - BD^{-1}CM = I \Rightarrow M = (A - BD^{-1}C)^{-1}$$

$$\text{And } P = - (D - CA^{-1}B)^{-1}CA^{-1}$$

In our case: $\Lambda^{-1} = \begin{bmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{bmatrix}^{-1} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$

$$\Sigma_{11} = (\Lambda_{11} - \Lambda_{12}\Lambda_{22}^{-1}\Lambda_{21})^{-1}$$

$$\Rightarrow Y_1 \sim N_{\mu_1}(\mu_1, \Sigma_{11})$$

Marginal density of Y_1

Main idea

- Marginal dist.
- Conditional dist

of
Multivariate
Normal

End of Week I

Week II

Conditional Density: \Rightarrow Multivariate Normal (from last week)

$$\underbrace{\left[\begin{array}{c} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{array} \right]}_{\mathbf{Y}} \quad p(\mathbf{Y}_2 | \mathbf{Y}_1) = \frac{p(\mathbf{Y}_1, \mathbf{Y}_2)}{p(\mathbf{Y}_1)}$$

Now, $\frac{p(\mathbf{Y}_1, \mathbf{Y}_2)}{p(\mathbf{Y}_1)} \propto p(\mathbf{Y}_2, \mathbf{Y}_1) \quad \text{w.r.t. } \mathbf{Y}_1$

$$\propto \exp\left(-\frac{1}{2} \delta_{\mathbf{Y}_2}^T \Lambda \delta_{\mathbf{Y}_2}\right) \quad \delta_{\mathbf{Y}_1} = \mathbf{Y}_1 - \boldsymbol{\mu}_1$$

$$\Rightarrow p(\mathbf{Y}_2 | \mathbf{Y}_1) \propto \exp\left(-\frac{1}{2} \left(\Sigma_{22} + \Lambda_{22}^{-1} \Lambda_{21} \Sigma_{12} \right)^T \Lambda_{22} \left(\Sigma_{22} + \Lambda_{22}^{-1} \Lambda_{21} \Sigma_{12} \right)\right)$$

$$= \exp\left[-\frac{1}{2} (\mathbf{y}_2 - \boldsymbol{\mu}_{2|1})^T \Lambda_{22} (\mathbf{y}_2 - \boldsymbol{\mu}_{2|1})\right]$$

where $\boldsymbol{\mu}_{2|1} = \boldsymbol{\mu}_2 - \Lambda_{22}^{-1} \Lambda_{21} (\mathbf{y}_1 - \boldsymbol{\mu}_1)$

From block inverse

\hookrightarrow But: 1) $\Lambda_{22} = (\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})^{-1}$

2) $\Lambda_{21} = -(\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})^{-1} \Sigma_{21} \Sigma_{11}^{-1}$

$\Rightarrow \boldsymbol{\mu}_{2|1} = \boldsymbol{\mu}_2 - \underbrace{(\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})}_{\Delta_{22}} \cdot \underbrace{(\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})^{-1}}_{\Delta_{21}}$

(Two relations given identity only one is used)

$$\underbrace{\Sigma_{21} \Sigma_{11}^{-1} (\mathbf{y}_1 - \boldsymbol{\mu}_1)}_{\Delta_{21}} = \underbrace{\Delta_{21}}_{\Delta_{21}}$$

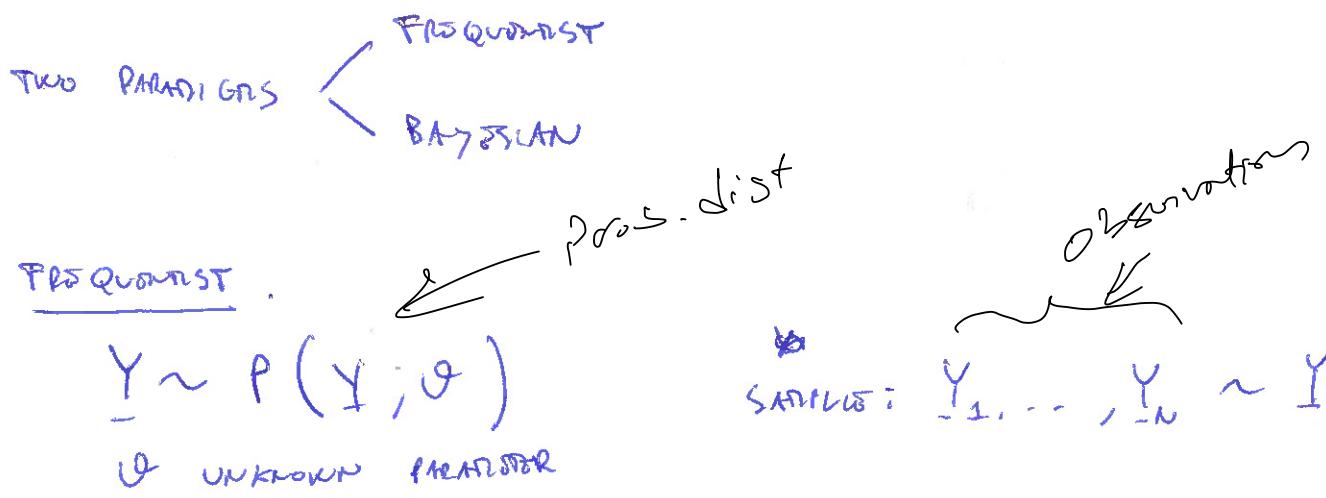
Conditional density $= \boldsymbol{\mu}_2 + \Sigma_{21} \Sigma_{11}^{-1} (\mathbf{y}_1 - \boldsymbol{\mu}_1)$

$$\Lambda_{22}^{-1}$$

$\Rightarrow \mathbf{y}_2 | \mathbf{y}_1 \sim N_p(\boldsymbol{\mu}_2 + \Sigma_{21} \Sigma_{11}^{-1} (\mathbf{y}_1 - \boldsymbol{\mu}_1), \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12})$

2 ESTIMATION

- MLE \rightarrow Frequentist
- Bayesian estimation for Gaussian



ESTIMATOR: FUNCTION OF THE SAMPLE $\hat{\theta} = g(Y_1, \dots, Y_n)$

(which hopefully is telling us something about θ)

We choose ^{the} estimator based on some (frequentist) properties:

- STABILITY $\text{TLSSE}(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \theta)^2] \leftarrow \text{MSE} \Rightarrow \text{Mean Sq. Err}$
- $\hat{\theta} \rightarrow \theta$ when $N \rightarrow \infty$

However, there are criterions to judge an estimator but they do not tell us how to pick one.

there are a few constructive procedure to build estimators with good ^(frequentist) properties, one of which is MAXIMUM LIKELIHOOD ESTIMATION

MLE in general:

$$\hat{\theta} = \arg \max L(\theta; Y_1, \dots, Y_N) \stackrel{\text{INDEPENDENT SAMPLE}}{\downarrow} = \arg \max \prod_{i=1}^N p(Y_i; \theta)$$

MLE for Gaussian Samples

\underline{Y}_l c.c.d. = independent and identically distributed

$$\underline{Y}_l \sim N_p(\underline{\mu}, \Sigma), \quad l=1, \dots, N$$

$$\theta = (\underline{\mu}, \Sigma)$$

- Likelihood = fⁿ of parameters given sample
- PDF = fⁿ of variables having certain parameters

Likelihood of Gaussian

$$\left\{ L(\underline{\mu}, \Sigma; \underline{Y}) = \prod_{l=1}^N \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{Y}_l - \underline{\mu})^\top \Sigma^{-1} (\underline{Y}_l - \underline{\mu})\right) \right.$$

$$\left. = \frac{1}{(2\pi)^{\frac{np}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left(-\frac{1}{2} \sum_{l=1}^N (\underline{Y}_l - \underline{\mu})^\top \Sigma^{-1} (\underline{Y}_l - \underline{\mu})\right) \right\}$$

To maximize it, it is useful to apply a log-transform:

$$L = \log(L) = -\frac{np}{2} \log(2\pi) - \frac{n}{2} \log(|\Sigma|) - \frac{1}{2} \sum_{l=1}^N (\underline{Y}_l - \underline{\mu})^\top \Sigma^{-1} (\underline{Y}_l - \underline{\mu})$$

Now, $(\underline{Y}_l - \underline{\mu})^\top \Sigma^{-1} (\underline{Y}_l - \underline{\mu}) = \underline{\mu}^\top \Sigma^{-1} \underline{Y}_l - 2 \underline{\mu}^\top \Sigma^{-1} \underline{\mu} + \underline{\mu}^\top \Sigma^{-1} \underline{\mu}$

To find the maximum: Differentiate log-likelihood & set to zero.

$$\boxed{\frac{\partial}{\partial \underline{\mu}} L = 0}$$

$$\Rightarrow \cancel{\text{Differentiation}}$$

$$-\frac{1}{2} \sum_{l=1}^N (-2 \Sigma^{-1} \underline{Y}_l + 2 \Sigma^{-1} \underline{\mu}) = 0$$

$$\Sigma^{-1} \left(\sum_{l=1}^N \underline{Y}_l - N \underline{\mu} \right) = 0$$

NOTE: IT IS A MAXIMUM
BECAUSE Σ IS A
POSITIVE DEFINITE QUADRATIC
FORM IN $\underline{\mu}$.

$\hat{\underline{\mu}} = \frac{1}{N} \sum_{l=1}^N \underline{Y}_l$

for Σ ,

$$\frac{\partial}{\partial \Sigma} \ell = 0 \Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\underline{Y}_n - \underline{\mu})^\top (\underline{Y}_n - \underline{\mu})$$

(left for exercise)

Also note: $\hat{\mu} \sim N_p(\mu, \frac{1}{N} \Sigma)$ (property of Gaussian distribution)

Bayesian Estimation:

You have prior knowledge, or belief, about the parameter θ that can be modelled with some probability distribution.

$$\theta \sim P(\theta)$$

then, you can use Bayes theorem to update your ~~old~~ belief about θ :

$$* \quad \underbrace{P(\theta | Y_1, \dots, Y_N)}_{\text{posterior}} \propto \underbrace{P(Y_1, \dots, Y_N | \theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}$$

and you get a posterior distribution for θ .

Bayesian estimation for Gaussian distributions

- Σ known (for simplicity):

Prior: $\underline{\mu} \sim N_p(\mu_0, \Sigma_0)$

Posterior: $P(\underline{\mu} | \underline{Y}) \propto P(\underline{Y} | \underline{\mu}) P(\underline{\mu}) = \dots$

↑ Using * above

$$-\propto \exp \left[-\frac{1}{2} \sum_{i=1}^N (\underline{Y}_i - \underline{Y}_0)^T \underline{\Sigma}^{-1} (\underline{Y}_i - \underline{Y}_0) \right] \exp \left[-\frac{1}{2} (\underline{Y} - \underline{Y}_0)^T \underline{\Sigma}_0^{-1} (\underline{Y} - \underline{Y}_0) \right]$$

$$= \exp \left[-\frac{1}{2} \left(\sum_{i=1}^N \underline{Y}_i^T \underline{\Sigma}^{-1} \underline{Y}_i - 2 \underbrace{\sum_{i=1}^N \underline{Y}_i^T \underline{\Sigma}^{-1} \underline{Y}} + \underbrace{N \underline{Y}^T \underline{\Sigma}^{-1} \underline{Y}} + \right. \right.$$

$$\left. \left. + \underline{Y}_0^T \underline{\Sigma}_0^{-1} \underline{Y} + \underline{Y}_0^T \underline{\Sigma}_0^{-1} \underline{Y}_0 - 2 \underline{Y}_0^T \underline{\Sigma}_0^{-1} \underline{Y} \right) \right]$$

(ignoring ~~term~~)

does not depend
on $\underline{\mu}$)

$$\propto \exp \left(-\frac{1}{2} \left[\underline{Y}^T \left(N \underline{\Sigma}^{-1} + \underline{\Sigma}_0^{-1} \right) \underline{Y} - 2 \left(\sum_{i=1}^N \underline{Y}_i^T \underline{\Sigma}^{-1} \underline{Y} + \underline{Y}_0^T \underline{\Sigma}_0^{-1} \underline{Y} \right) \right] \right)$$

$$\propto \exp \left(-\frac{1}{2} \left[\left(\underline{Y} - \left(\underline{\Sigma}_0^{-1} + N \underline{\Sigma}^{-1} \right)^{-1} \left(\sum_{i=1}^N \underline{Y}_i + N \underline{Y}_0 \right) \right)^T \cdot \right. \right.$$

Bayesian Estimation of $\underline{\mu}$ & $\underline{\Sigma}$
for Gaussian

$$\cdot \left. \left(N \underline{\Sigma}^{-1} + \underline{\Sigma}_0^{-1} \right) \left(\underline{Y} - \left(\underline{\Sigma}_0^{-1} + N \underline{\Sigma}^{-1} \right)^{-1} \left(\sum_{i=1}^N \underline{Y}_i + N \underline{Y}_0 \right) \right) \right]$$

↓

$$\Rightarrow \left\{ \begin{array}{l} \underline{Y} | Y \sim N_p(\underline{Y}_0, \underline{\Sigma}_0) \\ \underline{Y}_0 = \left(\underline{\Sigma}_0^{-1} + N \underline{\Sigma}^{-1} \right)^{-1} \left(\sum_{i=1}^N \underline{Y}_i + N \underline{Y}_0 \right) \\ \underline{\Sigma}_0 = \left(\underline{\Sigma}_0^{-1} + N \underline{\Sigma}^{-1} \right)^{-1} \end{array} \right\}$$

If you want a point estimator (as opposed to a distribution),

\underline{Y}_0 is the natural choice $\frac{\text{Mode}}{\text{Mode}}$. (because it is both the mean
and the mode of the posterior)

Note: $\underline{Y}_0 \rightarrow \underline{Y}_{\text{MLE}}$ if $N \rightarrow +\infty$ ← more data & less certainty of prior
Also: $\underline{Y}_0 \rightarrow \underline{Y}_0^{\text{prior}}$ if $N \rightarrow 0$ ← less data & more certainty of prior

- Gaussian priors are called CONJUGATE PRIOR Priors
- THE GAUSSIAN LIKELIHOOD (Posterior can be obtained with simple UPDATE RULE)

- IF Σ IS UNKNOWN,

PRIOR : $\underline{H} \sim N_p(H_0, \frac{1}{\lambda} \Sigma_{\text{pri})}$
 $\Sigma \sim \text{inverse wishart } W^{-1}(\gamma, v)$

Posterior : ~~$\underline{H} \sim N(\bar{H}, \frac{1}{\lambda} \Sigma)$~~
 $\Sigma | \underline{Y} \sim W^{-1}(\hat{\Sigma} + \bar{\Sigma}, N + v)$

$$\bar{H} = \frac{\lambda H_0 + N \bar{H}_{\text{ML}}$$

- In practice, conjugate prior are not used anymore.
 Computational methods are used to sample from posterior distributions;
 \hookrightarrow MARKOV CHAIN MONTE CARLO (MCMC)
- Prior which we do not have explicit expression.

- If I have more data & less certainty of prior
 then M^* will be the MLE.
 & If I have less of data & more certainty of prior then M^* will be closer to my prior.

Decision Theory

One way to know estimators from Posterior dist^t
One way to choose estimators from Posterior distributions is with
Decision theory:

Let \hat{Y}_t be an estimator of Y_t and assume a quadratic loss, i.e.

$$\text{L}(\hat{Y}_t) = (\hat{Y}_t - Y_t)^2 \quad \text{for } t=1, \dots, p$$

We would like to minimize ~~L~~ L , but Y_t is unknown. So we
minimize the posterior expected loss: (also called Bayes Risk)

$$\Rightarrow \mathbb{E}_{\hat{Y}_3|Y} [L(\hat{Y}_3)] = \mathbb{E}_{Y_3|Y} [L(\hat{Y}_3)] = \mathbb{E}_{Y_3|Y} [(\hat{Y}_3 - Y_3)^2] =$$

$$= \mathbb{E}_{Y_3|Y} [((\hat{Y}_3 - \mu_{\circledast 3}) + (Y_{\circledast 3} - Y_3))^2] = \quad \begin{matrix} \leftarrow \text{add & sub} \\ \text{max posterior } M^* \text{ at} \\ \text{before from Bayesian estimation} \end{matrix}$$

$$= \mathbb{E}_{Y_3|Y} [(\hat{Y}_3 - \mu_{\circledast 3})^2] + (\mu_{\circledast 3} - \hat{Y}_3)^2 + 2(\hat{Y}_3 - \mu_{\circledast 3}) \mathbb{E}[Y_3 - \mu_{\circledast 3}]$$

$$= V_{Y_3|Y}(\mu_{\circledast 3}) + (\mu_{\circledast 3} - \hat{Y}_3)^2$$

0, because
 $\mathbb{E}[Y_{\circledast 3}] = \mu$

This is minimized for $\hat{Y}_3 = \mu_{\circledast 3}$

minimizes

variance of posterior M_3

End of Week II

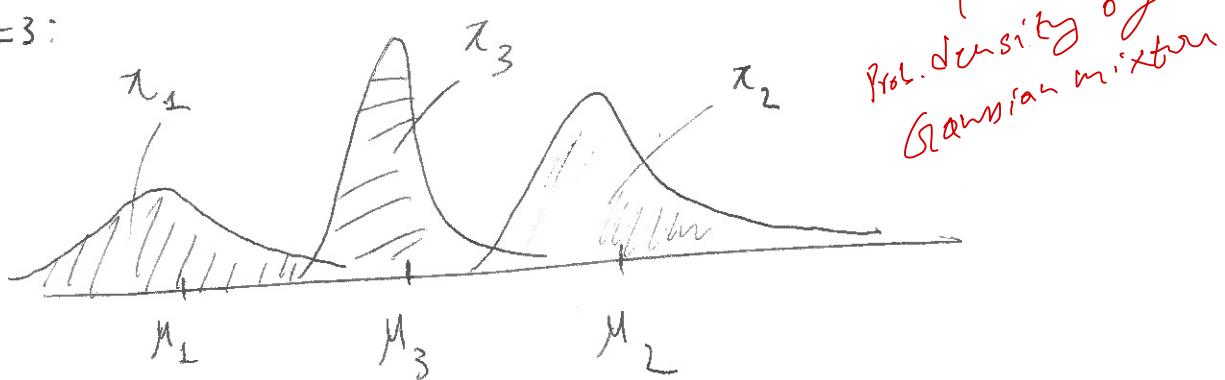
3. Gaussian Mixture Models

\hookrightarrow & Parameters estimation

Random variables that are "mixture of normals", i.e. they have probability distribution

$$P(\underline{y}; \theta) = \sum_{k=1}^K \pi_k \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left(-\frac{1}{2} (\underline{y} - \underline{\mu}_k)^T \Sigma_k^{-1} (\underline{y} - \underline{\mu}_k) \right)$$

In 1D; $K=3$:



$$\sum_{k=1}^K \pi_k = 1, \quad \theta = (\pi_1, \dots, \pi_K, \underline{\mu}_1, \dots, \underline{\mu}_K, \Sigma_1, \dots, \Sigma_K)$$

~~Value function~~

Alternative representation:

let \underline{w} be a ~~random~~ random vector such that:

$$w_k = \begin{cases} 1 & \text{if the observation is generated from the } k\text{-th component} \\ 0 & \text{else} \end{cases} \quad \text{comes from the } k\text{-th component}$$

i.e.: $P(\underline{w} = (0, \dots, \overset{k\text{-th}}{1}, \dots, 0)) = \pi_k$, or $[P(\underline{w}) = \prod_{k=1}^K \pi_k^{w_k}]$

$$\Rightarrow \underline{y} = \sum_{k=1}^K w_k \underline{y}_k, \quad \underline{y}_k \sim N(\underline{\mu}_k, \Sigma_k)$$

Conditional density of $\underline{Y} | \underline{w}$ is:

$$P(\underline{Y} | \underline{w}_k = 1) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{y} - \underline{\mu}_k)^\top \Sigma_k^{-1} (\underline{y} - \underline{\mu}_k)\right)$$

∴ $P(\underline{Y} | \underline{w}) = \prod_{k=1}^K P(\underline{Y} | w_k)$

And the Joint Density is:

$$P(\underline{Y}, \underline{w}) = \prod_{k=1}^K \pi_k \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (\underline{y} - \underline{\mu}_k)^\top \Sigma_k^{-1} (\underline{y} - \underline{\mu}_k)\right)$$

Why are these useful?

- .) Flexible model to represent data that are multimodal
- .) Easy to simulate from: Pick one component with prob π_k and then sample from $N(\underline{\mu}_k, \Sigma_k)$

.) CLUSTERING (UNSUPERVISED LEARNING)

Knowing that a data point is likely from one component is a way to cluster (i.e. group) the data.

$$P(w_k | \underline{Y}) = \frac{P(\underline{Y}, w_k)}{P(\underline{Y})} = \frac{P(w_k) P(\underline{Y} | w_k)}{\sum_{i=1}^K P(w_i) P(\underline{Y} | w_i)} =$$

$$= \frac{\pi_k P(\underline{Y} | w_k = 1)}{\sum_{i=1}^K \pi_i P(\underline{Y} | w_i = 1)}$$

⇒ LABEL $\hat{k} = \arg \max_{k=1, \dots, K} P(w_k | \underline{Y})$

To do this, I need to know the parameters $\underline{\theta}$; this is where we apply the parameter estimation techniques we have seen last time.

Let $\underline{Y}_1, \dots, \underline{Y}_N$ be i.i.d. with the same density as \underline{Y} . (Gaussian mixtures with parameters $\underline{\theta}$). Then, we have

Latent variables (also called hidden variables)

$$w_1, \dots, w_N \sim \underline{w}$$

They represent the component of the mixtures from which each data point comes from. Unfortunately, we do not ~~observe~~ observe them.

(Aim): Find the MLE for $\underline{\theta}$

$$L(\underline{\theta}; \underline{Y}, \underline{w}) = \prod_{i=1}^N \prod_{k=1}^K \left\{ \frac{\pi_k}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\underline{y}_i - \underline{\mu}_k)^T \Sigma^{-1} (\underline{y}_i - \underline{\mu}_k) \right] \right\}$$

$$\Rightarrow \ell = \log(L) = \sum_{i=1}^N \left\{ \sum_{k=1}^K w_{ik} \left[\log \pi_k - \frac{1}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) + \right. \right. \\ \left. \left. - \frac{1}{2} (\underline{y}_i - \underline{\mu}_k)^T \Sigma^{-1} (\underline{y}_i - \underline{\mu}_k) \right] \right\}$$

This would be easy to maximize if we knew w_1, \dots, w_K .

But we don't.

↳ By separating data coming from each $\xrightarrow{\text{each component}} \xrightarrow{\text{separately}}$ component and estimating the mixtures of other components separately

$$\oplus \quad \hat{\pi}_k = \frac{1}{N} \sum_{i=1}^N w_{ik} \quad \text{estimating}$$

⇒ Solve the problem using Expectation-Maximization ALGORITHM

Parenthesis : Information Theory And KULLBACK-LEIBLER (KL)

Divergence

Information Theory studies the transmission of information. the amount of information given by \underline{Y} is measured by the entropy:

$$H(\underline{Y}) = - \int p(y) \lg(p(y)) dy$$

If we have a second distribution $q(y)$, which is used to approximate $p(y)$, we can define the relative entropy or KL Divergence between p and q : (i.e. Answer of question lost when q is used to approximate p)

$$KL(p \parallel q) = - \int p(y) \lg(q(y)) dy + \int p(y) \lg(p(y)) dy$$

$$KL = - \int p(y) \lg \left[\frac{q(y)}{p(y)} \right] dy$$

Note that: $KL(q \parallel p) \neq KL(p \parallel q)$

- By Jenson INEQUALITY

$$\int p(y) \lg \left[\frac{q(y)}{p(y)} \right] \leq \lg \underbrace{\int p(y) \frac{q(y)}{p(y)} dy}_{1} = 0$$

$$\Rightarrow \boxed{KL(p \parallel q) \geq 0}$$

(First part we are going to use.)

- $KL(p \parallel q) = 0 \Leftrightarrow p(y) = q(y)$ almost everywhere

EM ALGORITHM

$$\text{Aim: } \max_{\underline{\theta}} \left[\log L(\underline{\theta}; \underline{Y}) \right] = \max_{\underline{\theta}} p(Y | \underline{\theta})$$

$$\text{where } p(Y | \underline{\theta}) = \sum_{\underline{w}} p(Y, \underline{w} | \underline{\theta})$$

HARD TO MAXIMIZE

easy, but \underline{w} is "HIDDEN"

~~Repetitive calculations~~

General IDEA of the ALGORITHM:

0.) START with some value of $\underline{\theta}$

1.) ESTIMATING THE PROBABILITY DISTRIBUTION OF \underline{w} , based on the current value of $\underline{\theta}$ (E-step)

2.) ESTIMATING A NEW VALUE OF $\underline{\theta}$, by maximizing

$$\sum_{\underline{w}} p(Y | \underline{\theta}, \underline{w}) p(\underline{w}) \quad (\text{M-step})$$

3.) REPEAT 1-2 UNTIL CONVERGENCE

More formally, we introduce a "trial" distribution $q(\underline{w})$ and we take advantage of the following identity:

$$\log p(Y|\theta) = \log \left(\frac{p(Y, \underline{w}|\theta)}{p(\underline{w}|Y, \theta)} \right) = \underline{Y} \cdot \underline{w}$$

$$= \log \left(\frac{p(Y, \underline{w}|\theta)}{q(\underline{w})} \frac{q(\underline{w})}{p(\underline{w}|Y, \theta)} \right) =$$

$$= \underbrace{\log \left(\frac{p(Y, \underline{w}|\theta)}{q(\underline{w})} \right)}_{\text{Expectation}} + \log \left(\frac{q(\underline{w})}{p(\underline{w}|Y, \theta)} \right)$$

If we take the expectation w.r.t. $q(\underline{w})$:

$$\log p(Y|\theta) = \underbrace{\sum_w q(w) \log \left[\frac{p(Y, \underline{w}|\theta)}{q(w)} \right]}_+ +$$

$$Q(\underline{q}, \underline{\theta}) = \underbrace{- \sum_w q(w) \log \left[\frac{p(\underline{w}|Y, \theta)}{q(w)} \right]}_-$$

$$\Rightarrow \boxed{\log p(Y|\theta) > Q(\underline{q}, \underline{\theta})} \quad KL(q||p) \geq 0$$

BTW iterations: start $\underline{q}^{(\text{old})}, \underline{\theta}^{(\text{old})}$

$$\max_q Q(q, \underline{\theta}^{(\text{old})}) = q^{\text{new}} \quad (\text{S-step})$$

$$\max_{\underline{\theta}} Q(q^{\text{new}}, \underline{\theta}) = \underline{\theta}^{(\text{new})} \quad (\text{R-step})$$

$$\text{E-step: } Q(q, \vartheta^{(\text{old})}) = \log \left\{ p(Y | \vartheta^{(\text{old})}) \right\} - KL(q || p^{(\text{old})})$$

this is maximized when $KL(q || p^{(\text{old})}) = 0$

$$\Rightarrow q^{\text{new}} = p(x | Y, \vartheta^{(\text{old})})$$

$$\text{M-step: } \underset{\vartheta}{\operatorname{Max}} \quad Q(q^{\text{new}}, \vartheta)$$

Note that $Q(q^{\text{new}}, \vartheta^{\text{new}}) \geq Q(q^{\text{new}}, \vartheta^{(\text{old})})$,

$$\begin{aligned} \text{Then: } \log p(Y | \vartheta^{\text{new}}) &= \underbrace{Q(q^{\text{new}}, \vartheta^{\text{new}})}_{\geq Q(q^{\text{new}}, \vartheta^{(\text{old})})} + \underbrace{KL(q^{\text{new}} || p^{\text{new}})}_{\geq 0} \\ &\geq Q(q^{\text{new}}, \vartheta^{(\text{old})}) \end{aligned}$$

$$\text{and } \log \left\{ p(Y | \vartheta^{(\text{old})}) \right\} = Q(q^{\text{new}}, \vartheta^{(\text{old})}) + \underbrace{KL(q^{\text{new}} || p^{(\text{old})})}_{= 0}$$

$$\Rightarrow \log \left\{ p(Y | \vartheta^{\text{new}}) \right\} \geq \log \left\{ p(Y | \vartheta^{(\text{old})}) \right\}$$

i.e. each iteration of the E-M algorithm is guaranteed to increase the log-likelihood, i.e. it converges to a local maximum of $\log p(Y | \vartheta)$

E-step:

$$\textcircled{B} \quad Q(q^{\text{new}}, \theta) = \sum_w q^{\text{new}}(w) \log \left[\frac{p(y, w | \theta)}{q^{\text{new}}(w)} \right] =$$

\star

$$= \underbrace{\sum_w q^{\text{new}}(w) \log p(y, w | \theta)}_{\text{const}} + \text{const}$$

→ $\tilde{Q}(\underline{\theta}, \underline{\theta^{\text{old}}})$

EM for (Gaussian) mixtures:

- $\log p(y, w | \theta) = \sum_{l=1}^N \sum_{k=1}^K w_{lk} \left[\log \pi_k - \frac{1}{2} \log |\Sigma_k| + - \frac{1}{2} (\underline{y}_l - \underline{\mu}_k)^T \Sigma_k^{-1} (\underline{y}_l - \underline{\mu}_k) \right]$

this following is from pg 2.

- $p(w | y, \theta) = \prod_{l=1}^N p(w_l | \underline{y}_l, \theta) \sim N(\underline{\mu}_k, \Sigma_k)$

AND $p(w_{lk}=1 | \underline{y}_l, \theta) = \frac{\pi_k p(y | w_k=1)}{\sum_{l=1}^K \pi_l p(y | w_l=1)} = \gamma_k(\underline{y}_l)$

$\gamma_k(\underline{y}_l)$ is called responsibility of component k for the data point \underline{y}_l .

γ^{new} from E-step

responsibility

from M-step last page (④) & Suss. • & oo last page

$$\Rightarrow \tilde{Q}(\boldsymbol{\theta}, \boldsymbol{\vartheta}^{old}) = \sum_{\mathbf{w}} p(\mathbf{w} | \mathbf{Y}, \boldsymbol{\vartheta}^{old}) \log p(\mathbf{y}, \mathbf{w} | \boldsymbol{\vartheta}^{old}) =$$
$$= \sum_{l=1}^N \sum_{k=1}^K \gamma_k(y_l) \left[\log \pi_k - \frac{1}{2} \log |\Sigma_k| + \right. \\ \left. - \frac{1}{2} (y_l - \mu_k)^T \Sigma_k^{-1} (y_l - \mu_k) \right]$$

We can maximize separately in π_k and in μ_k, Σ_k :

$$\max_{\mu_k, \Sigma_k} \sum_{l=1}^N \gamma_k(y_l) \log \left\{ \frac{1}{(2\pi)^{k/2} |\Sigma_k|^{k/2}} \exp \left\{ -\frac{1}{2} (y_l - \mu_k)^T \Sigma_k^{-1} (y_l - \mu_k) \right\} \right\}$$

Maximun is $\frac{N}{N_k}$:

$$\mu_k^{new} = \frac{1}{N_k} \sum_{l=1}^N \gamma_k(y_l) y_l, \quad N_k = \sum_{l=1}^N \gamma_k(y_l)$$

Max. with Σ :

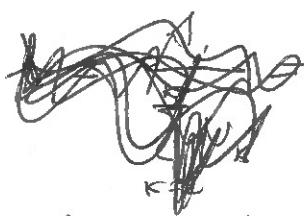
$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{l=1}^N \gamma_k(y_l) (y_l - \mu_k) (y_l - \mu_k)^T$$

$$\max \sum_{l=1}^N \sum_{k=1}^K \gamma_k(y_l) \log \pi_k = \sum_{k=1}^K N_k \log \pi_k$$

with $\sum_{k=1}^K \pi_k = 1$. Using a lagrange multiplier:

$$\max_{\pi} \left[\sum_{k=1}^K N_k \log(\pi_k) - \lambda \left(\sum_{k=1}^K \pi_k - 1 \right) \right] \Rightarrow$$

$$\Rightarrow \frac{\partial}{\partial \lambda_k} (*) = \frac{N_k}{\lambda_k} - \lambda = 0 \Rightarrow \lambda_k = \frac{N_k}{\lambda}$$



$$\sum_{k=1}^K \lambda_k = 1 \Rightarrow \sum_{k=1}^K \frac{N_k}{\lambda} = 1 \Rightarrow \lambda = \sum_{k=1}^K N_k = \sum_{k=1}^K \sum_{i=1}^n f_k(y_i) = N$$

\Rightarrow

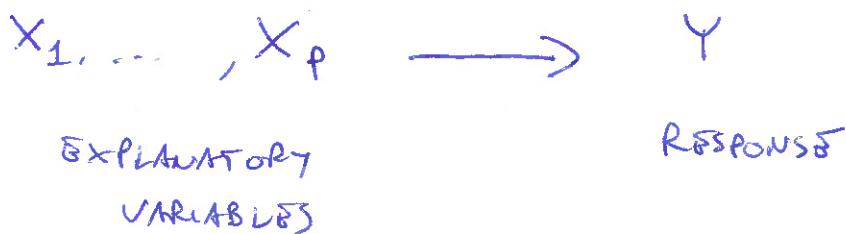
$$\lambda_k^{\text{new}} = \frac{N_k}{N}$$

End of Week III

4.) LINEAR MODELS

Week IV

We start addressing the problem of prediction:



From now
it will
be
Prediction
Problem.

We want to predict the response Y from the explanatory variables, i.e. $Y = f(x_1, \dots, x_p)$.

To LEARN the relationship f , we use a SAMPLE.

$$(x_{1c}, \dots, x_{nc}, y_c) \quad , \quad c=1, \dots, N$$

GAUSSIAN LISSOMA MODEL

ϕ_i known,
 β_i unknown.

$$Y_v | X_v = x_v \stackrel{\text{L.L.d.}}{\sim} N(\mu_v, \sigma^2), \quad v=1, \dots, N$$

where $M_v = \beta_1 \phi_1(x_v) + \beta_2 \phi_2(x_v) + \dots + \beta_k \phi_k(x_v)$

MATRIX FORM

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_p \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_p \end{bmatrix}$$

$$X = \begin{bmatrix} \phi_1(x_1) & \dots & \phi_k(x_1) \\ \vdots & & \vdots \\ \phi_1(x_N) & \dots & \phi_k(x_N) \end{bmatrix}$$

$$\hat{Y} = X\beta + \varepsilon ; \quad \varepsilon \sim N_n(0, \sigma^2 I), \varepsilon \text{ indep.}$$

γ ~ $N_p(X\beta, \sigma^2 I)$ \rightarrow errors

Remark: In OLS geometry, we can use a linear model

$$\underline{Y} = \underline{X}\beta + \underline{\epsilon}$$

where $\underline{\epsilon}$ are independent but not necessarily Gaussian, most of the results will still hold.

LEAST SQUARES ESTIMATION

We can estimate β by minimizing the residual ~~sum~~^{sum} of squares:

$$S(\beta) = \sum_{i=1}^N \left(Y_i - \sum_{j=1}^p \phi_j(x_i) \beta_j \right)^2 = (\underline{Y} - \underline{X}\beta)^T (\underline{Y} - \underline{X}\beta)$$

$$\frac{\partial S}{\partial \beta_k} = \frac{\partial}{\partial \beta_k} \sum_{i=1}^N \left[Y_i - \sum_{j=1}^p \phi_j(x_i) \beta_j \right]^2 =$$

$$= -2 \sum_{i=1}^N \left[Y_i - \sum_{j=1}^p \phi_j(x_i) \beta_j \right] \cdot \phi_k(x_i)$$

$$\frac{\partial S}{\partial \beta} = -2 \underline{X}^T (\underline{Y} - \underline{X}\beta) = 0$$

$$\underline{X}^T (\underline{Y} - \underline{X}\beta) = 0$$

$$\hat{\beta} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$

IF $(\underline{X}^T \underline{X})$ is INVERTIBLE

IF $\underline{X}^T \underline{X}$ is NOT INVERTIBLE \rightarrow NOT UNIQUE SOLUTION.

We will see how to deal with ~~this~~ this case ~~with~~ with REGULARIZED METHOD

Remark: If $\varepsilon \sim N_N(0, \sigma^2 I)$, the DLS for β is or the same as the least squares estimator:

$$L(\beta, \sigma^2; y) = \frac{1}{(2\pi)^{n/2} |\sigma^2 I|^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta) \right\}$$

$$\log L = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\sigma^2 I| - \frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta)$$

$$\begin{aligned} \underset{\beta}{\operatorname{argmax}} \log L &= \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta) = \\ &= (X^T X)^{-1} X^T Y \end{aligned}$$

- $\mathbb{E}[\hat{\beta}] = (X^T X)^{-1} X^T \mathbb{E}[Y] = (X^T X)^{-1} X^T X \beta = \beta$ UNBIASED ESTIMATOR

- $\text{Var}(\hat{\beta}) = (X^T X)^{-1} X^T \text{Var}(Y) X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-1} = \sigma^2 (X^T X)^{-1}$

$$\Rightarrow \hat{\beta} \sim N_p(\beta, \sigma^2 (X^T X)^{-1}) \quad (\text{since linear combination of Gaussian variables})$$

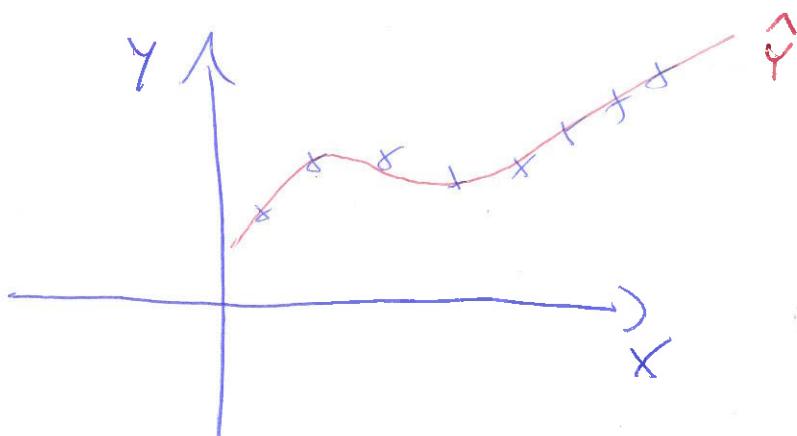
This DLS for σ^2 is:

$$\hat{\sigma}^2 = \frac{1}{N} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$$

(left for exercise) but we know this is a BIASED estimator.

UNBIASED VERSION : $s^2 = \frac{1}{N-p} (Y - X\hat{\beta})^T (Y - X\hat{\beta})$

EXAMPLE : POLYNOMIAL REGRESSION:



$$Y_i = \mu_i + \varepsilon_i$$

$$\mu_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3$$

$$\phi_0(x) = 1 \quad \phi_1(x) = x \quad \phi_2(x) = x^2 \quad \phi_3(x) = x^3$$

$$X = \begin{bmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & \vdots & \vdots & \vdots \\ 1 & x_N & x_N^2 & x_N^3 \end{bmatrix}$$

PREDICTED VALUES: $\hat{Y} = \hat{\beta}_0 \phi_0(x) + \hat{\beta}_1 \phi_1(x) + \hat{\beta}_2 \phi_2(x) + \hat{\beta}_3 \phi_3(x)$

EXTENSIONS OF LINEAR MODELS (two special cases)

- MULTIVARIATE RESPONSE: (when response \mathbf{Y} is vector itself)

$$\underline{Y}_i \stackrel{\text{iid}}{\sim} N_m(\underline{x}_i^T \underline{B}, \Sigma)$$

$$\underline{Y} = \mathbf{X} \underline{B} + \underline{\varepsilon} \leftarrow \text{error}$$

$$\underline{\varepsilon} \sim \text{Matrix Valued Normal } (\mathbf{0}_{nm}, \Sigma \otimes I_n)$$

LEAST SQUARES

$$\hat{\underline{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{Y}$$

(same as many linear regressions)

Alternatively,

~~WEIGHTED LEAST~~ GENERALIZED LEAST SQUARED

$$\hat{\underline{B}} = \underset{\Sigma}{\text{argmin}} \sum_{i=1}^n (\underline{y}_i - \underline{x}_i^T \underline{B})^T \Sigma^{-1} (\underline{y}_i - \underline{x}_i^T \underline{B})$$

PROBLEM: Σ needs to be estimated

Possible Solutions: ITERATIVE PROCEDURE

$\cdot \underline{B}^{(0)}, \Sigma^{(0)}$ — start value

~~WEIGHTED LEAST~~ $\cdot \Sigma^{(0)} \rightarrow \underline{B}^{(1+1)} — \text{estimate } \underline{B}^{(1+1)} \text{ using } \Sigma^{(0)}$

$\cdot \underline{B}^{(1+1)} \rightarrow \Sigma^{(1+1)} — \text{estimate } \Sigma^{(1+1)} \text{ using } \underline{B}^{(1+1)}$

- Correlated Response:

$$\underline{Y} \sim N_N(\mathbf{X} \underline{B}, \sigma^2 \Sigma)$$

IF Σ known, GENERALIZED LEAST SQUARED SOLUTION:

$$\hat{\underline{B}} = \underset{\Sigma}{\text{argmin}} (\underline{Y} - \mathbf{X} \underline{B})^T \Sigma^{-1} (\underline{Y} - \mathbf{X} \underline{B}) \Rightarrow \hat{\underline{B}} = (\mathbf{X}^T \Sigma^{-1} \mathbf{X})^{-1} \mathbf{X}^T \Sigma^{-1} \underline{Y}$$

• IF τ UNKNOWN, AS ABOVE we can use an iterative procedure.

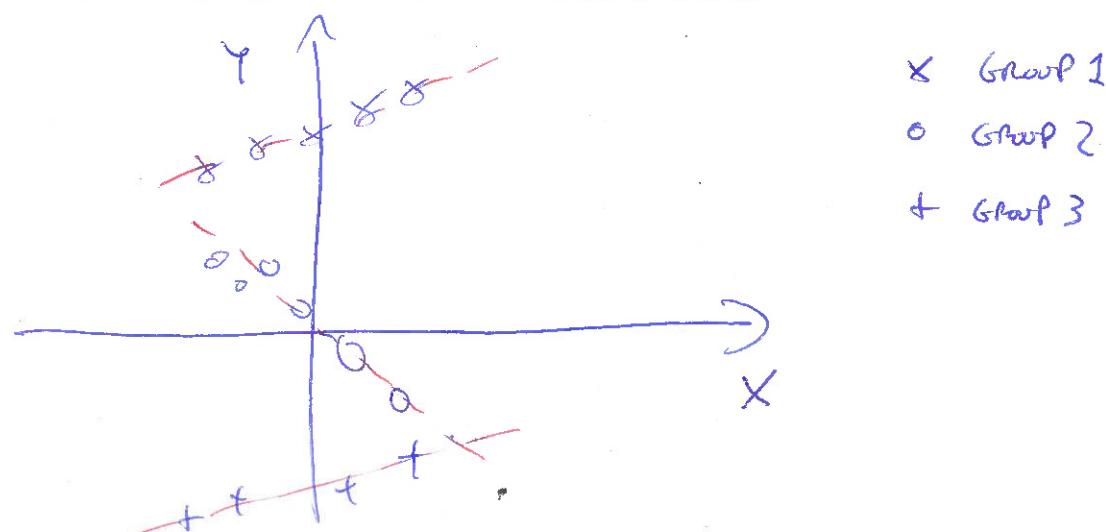
(POPULAR IN TIME SERIES AND SPATIAL STATISTICS)

but we need some assumptions on the structure of τ .

• (We cannot estimate a N dimensional covariance with 1 observation)

EXAMPLE: TIME SERIES : $\text{cov}(\varepsilon_i, \varepsilon_j) = \delta(1-i-j)$

• INDICATOR VARIABLES VS RANDOM EFFECTS



$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 I_{1i} + \beta_3 I_{2i} \quad \text{(crossed out)}$$

$$+ \beta_4 I_{1i} X + \beta_5 I_{2i} X + \varepsilon_i$$

INDICATOR
VARIABLES

$$I_{1i} = \begin{cases} 1 & \text{if } i \in \text{Group 2} \\ 0 & \text{else} \end{cases}$$

$$I_{2i} = \begin{cases} 1 & \text{if } i \in \text{Group 3} \\ 0 & \text{else} \end{cases}$$

⇒ THIS Allows TO FIT INDIVIDUAL LINES FOR EACH GROUP.

However, IF THE GROUPS ARE ESTIMATED FROM A LARGER POPULATION
(EXAMPLES: SCHOOLS, HOSPITALS, COUNTIES, ETC..), IT IS
BETTER TO USE RANDOM EFFECTS INSTEAD.

Linear Mixed Effects Models (Random effect)

$$\begin{cases} Y|u=\bar{u} \sim N(X\beta + Z\bar{u}, \sigma^2 I_n) \\ u \sim N(0, G) \end{cases} \quad \begin{matrix} \sigma \text{ is unknown} \\ \beta \text{ is unknown} \end{matrix}$$

\hat{y}

Random effect

$$Y = X\beta + Zu + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2 I_n), u \sim N(0, G)$$

$$\Rightarrow Y \sim N(X\beta, \Sigma), \quad \Sigma = ZGZ^\top + \sigma^2 I_n$$

$$\Rightarrow \hat{\beta} = (X\Sigma^{-1}X)^\top X^\top \Sigma^{-1} Y \quad \leftarrow \begin{matrix} \text{using} \\ (\text{Generalised}) \text{ Least Squares} \end{matrix}$$

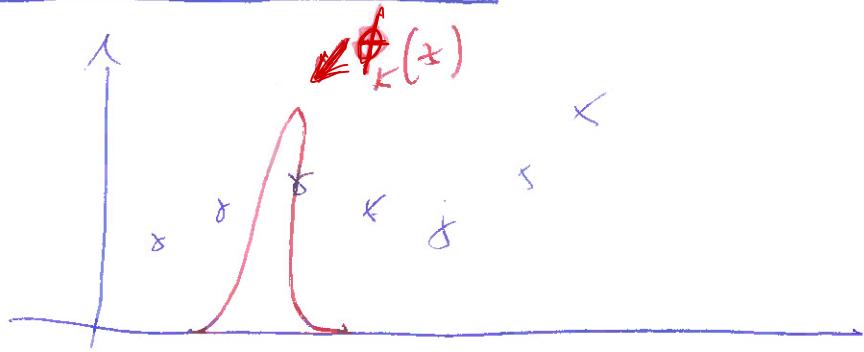
$$\hat{\Sigma} = \frac{\text{numerically}}{\text{numerically}} \quad \leftarrow \begin{matrix} \text{REML Algorithm} \\ \downarrow \\ \text{Restricted Max. likelihood Algorithm} \end{matrix}$$

(REML)

End of Week III

Week V

Local Basis Functions



(last week material continues)

$$Y_l = \sum_{j=1}^P f_j \phi_j(x) + \varepsilon_l$$

EXAMPLES: SPLINES, WAVELETS, LOCAL POLYNOMIAL

Additive Models

$$Y_l = \sum_{j=1}^P f_j(X_j) + \varepsilon_l \quad , l=1, \dots, n$$

↳ TO BE ESTIMATED SOMEHOW, FOR EXAMPLE:

$$f_j(x_j) = \sum_{n=1}^R \beta_{jn} \phi_n(x_j)$$

$$\Rightarrow Y_l = \sum_{j=1}^P \sum_{n=1}^R \beta_{jn} \phi_n(x_j) + \varepsilon_l$$

IN PRACTICE, FITTED VIA BACKFITTING ITERATIVE PROCEDURES; OR
 REGULARIZED METHODS
 Regularizes

5) MODEL SELECTION AND REGULARIZATION

(New material
for this week)

How do we choose the model to fit?

How to choose which and how many functions $\phi_s(\underline{x})$ to use?

Remark : Whenever possible, scientific knowledge about the variables of interest needs to be used to restrict as much as possible the set of models to be considered.

However, we will focus on how to use the data to select among a set of models.

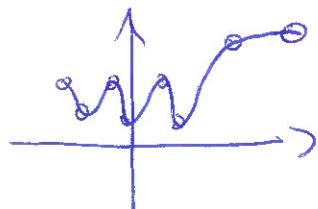
Example : Polynomial Regression: $\underline{Y} \sim N_n(\underline{\mu}, \sigma^2 \underline{I}_n)$

$$\underline{\mu}_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_i^p$$

How should we choose ~~the~~ p (degree of the polynomial)?

Note: $p=N-1 \Rightarrow s(f)=0$ (INTERPOLATING POLYNOMIAL)

However, the model is useless because ~~as~~ it is adapted to the error in our data \Rightarrow EXAMPLE OF OVERTFITTING:



TWO ASPECTS OF MODEL SELECTION:

- 1) • CRITERION TO COMPARE TWO MODELS
- 2) • SEARCH OF THE MODEL SPACE

For what concerns 2), if the model space is small you can compute your criterion for all models and pick the best one. If the model space is large, this can be not computationally feasible. A popular alternative is to use stepwise algorithm, which add or remove one predictor at the time from the model and they stop when the chosen criterion cannot be improved by adding or removing any term. It is a greedy approach, no guarantees to find the best model overall.

About 1), we are going to see a gallery of used criteria:

(i) - SPLIT DATA IN TRAINING / TEST SET

You FIT your models using the Training set and the criterion is the prediction error on the data in the Test set.

For example, $i=1, \dots, n \rightarrow$ TRAINING SET

$i=n+1, \dots, N \rightarrow$ TEST SET

use $y_1, \dots, y_n \rightarrow$ estimate \hat{p}

Choose the model which minimizes $\sum_{i=n+1}^N (y_i - x_i^T \hat{p})^2$

DRAWBACK: it wastes some data, less accurate estimates

• THE ESTIMATE OF THE PREDICTION ERROR IS ALSO ~~NOT~~
HIGHLY VARIABLE (SMALL TEST SAMPLE)

(ii) - CROSS-VALIDATION:

You repeat the ~~fold~~ partition between training and testing set multiple times and look at the ~~average~~ average prediction error

Special cases:

(a) K-Fold cross validation:

- Split the data in K sets

. for $k=1:K$:

. fit the model without the data in the k-th set

$$\text{compute } e_k = \frac{1}{|N_k|} \sum_{j \in N_k} (y_j - \hat{y}_j^{(-k)})^2$$

where N_k is the set of indices in the k-th partition

and $\hat{y}_j^{(-k)}$ is the prediction of the model fitted without the k-th set of data.

. choose the model that minimizes $\frac{1}{K} \sum_{k=1}^K e_k$ → average of predicted errors

(b) Leave-1-out CV:

. for $i=1:N$

. fit the model without the i-th observation

~~compute $e_i = \frac{1}{N} (y_i - \hat{y}_i^{(-i)})^2$~~

. choose the model that minimizes $\frac{1}{N} \sum_{i=1}^N e_i^2$

ADVANTAGES: You extract information from all the data

DRAWBACKS: - COMPUTATIONALLY EXPENSIVE

- YOU CAN STILL OVERFIT YOUR MODEL TO YOUR SPECIFIC DATASET.

(iii) Akaike Information Criterion (AIC)

This is one example of criteria that try to balance the prediction error and the complexity of the model:

$$AIC = 2p - 2 \log L(\hat{\theta}, \hat{\sigma}^2; Y)$$

\square p is the number of parameters in the model

Other information criteria use different penalties:

$$BIC: \log(N)P$$

$$DIC: \dots$$

- For large N , AIC:
 - approximate the leave-1-out CV error
 - \square the difference in AIC between two models is an estimate of the KL divergence between these models and the true model.
- choose model that minimizes AIC.

REGULARIZATION

Instead of selecting a model, we can use the largest possible model but shrinking the estimated coefficients to control their variance (and thus avoid overfitting).

RIDGE REGRESSION :

$$\hat{\beta}_R = \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta) =$$

$$\beta^T \beta \leq C \quad \xrightarrow{\text{SHRINKING CONSTRAINT}}$$

LAGRANGIAN MULTIPLIER

$$= \underset{\beta}{\operatorname{argmin}} \left\{ (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta \right\}$$

Now minimize $\tilde{S}(\beta; \lambda)$,

$$\frac{\partial \tilde{S}}{\partial \beta} = -2X^T Y + 2X^T X\beta + 2\lambda\beta = 0$$

$$\Rightarrow (X^T X + \lambda I)\beta = X^T Y$$

$$\boxed{\hat{\beta}_R = (X^T X + \lambda I)^{-1} X^T Y}$$

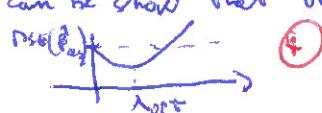
Remarks :

- Above, $C = \hat{\beta}_R^T \hat{\beta}_R$; λ is a tuning parameter to be chosen.

- $\hat{\beta}_R$ is well defined even if $X^T X$ is not invertible.

- $\lambda = 0 \Rightarrow \hat{\beta}_R = \hat{\beta}_{OLS}$; $\lambda \rightarrow +\infty \Rightarrow \hat{\beta}_R \rightarrow 0$

- It can be shown that the root sum of error as function of λ follows:



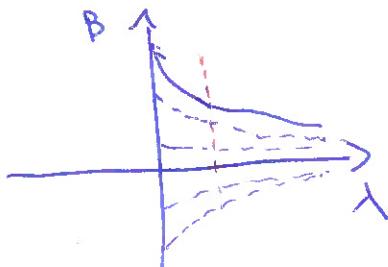
Remarks:

- The shrinking happens equally on all coefficients, so we need to standardize the variables ~~first~~ first.

- How to choose λ :

(i) CROSS-VARIATION

(ii) RIDGE TRACE



minimum λ that "stabilizes" the estimates

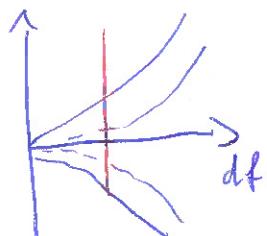
(iii) Approximation of λ_{opt} given some estimate of σ^2

Once λ is chosen then check:

- EFFECTIVE DEGREES OF FREEDOM: (Effective degrees of freedom of model.)

$$\hookrightarrow df(\lambda) = \text{trace} \left(X(X^T X + \lambda I)^{-1} X^T \right)$$

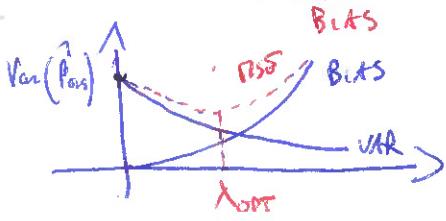
since $\hat{Y} = X(X^T X + \lambda I)^{-1} X^T Y$, it is a measure of
the dimension of the space to which fitted values belong.



If $\lambda \rightarrow \infty$, $df(\lambda) \rightarrow 0$
 $\lambda \rightarrow 0$, $df(\lambda) \rightarrow$ original no. of parameters.

* ASIDE: BIAS-VARIANCE TRADE-OFF

$$\begin{aligned} \text{MSE}(\hat{\beta}_R) &= \| E[\hat{\beta}_R - \beta] \|^2 + \text{trace}(\text{Var}(\hat{\beta}_R)) = \\ &= \lambda^2 \beta^T W^2 \beta + \sigma^2 \text{trace}(W(X^T X) W) \end{aligned}$$



$$\lambda = 0: \text{BIAS} = 0, \text{Var} = \frac{\text{Var}(\beta_{\text{OLS}})}{(\text{Var} + \infty)}$$

$$\lambda \rightarrow \infty: \text{BIAS} \rightarrow \infty, \text{Var} \rightarrow 0$$

(Engle-Watson II)

$$W = (X^T X + \lambda I)^{-1}$$

Week VII

(last week material
continued)

LASSO:

There is no particular reason to shrink the coefficients based on the 2-norm. The LASSO estimator is used.

$$\hat{\beta}_L = \underset{\beta}{\operatorname{arg\,min}} \quad (Y - X\beta)^T (Y - X\beta) =$$

$$\sum_i |\beta_i| \leq C \quad \leftarrow \text{constraint} \quad \|P\|_1$$

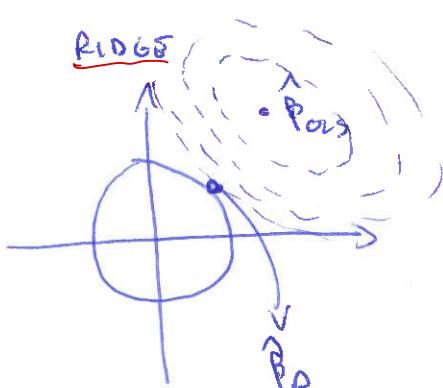
$$= \underset{\beta}{\operatorname{arg\,min}} \left\{ (Y - X\beta)^T (Y - X\beta) + \lambda \sum_i |\beta_i| \right\}$$

No explicit solution, it needs to be solved numerically
 (LARS algorithm) (Also called Norm 1 (lasso))

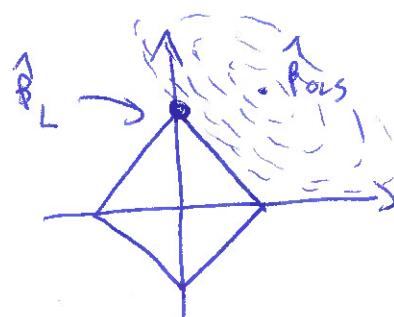
Drawbacks: • Again, variables need to be standardized first.

• $\lambda \rightarrow 0$, $\hat{\beta}_L \rightarrow \hat{\beta}_{OLS}$; $\lambda \rightarrow \infty$, $\hat{\beta}_L \rightarrow 0$

• VISUAL COMPARISON IN 2D:



LASSO



Minimum occurs at boundary of ellipse & circle or

ellipses & line.

LASSO ESTIMATOR favors SPARSITY.
favors

- Choice of λ based on similar strategies as in the RIDGE regression.
- In statistical software, often reparameterized as $s = \frac{\lambda}{\|\hat{\beta}_{OLS}\|_1}$, so that $0 < s < 1$.

- It is possible to play around further with the penalty, tailoring it to the problem and basis function used.

Example: Smoothing splines :

$$\hat{f} = \sum_{i=1}^n \hat{\beta}_i f_i(x)$$

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 + \lambda \int_a^b [\hat{f}''(t)]^2 dt$$

- Lasso is equivalent to soft-thresholding

$$\hat{\beta}_{i,j} = \begin{cases} 1 & \text{if } \beta_{i,j} > 0 \\ 0 & \text{if } \beta_{i,j} \leq 0 \end{cases}$$

IF THE PREDICTORS ARE ORTHOGONAL
(EXAMPLE: WAVELETS OR FOURIER BASIS)

- Best subset selection:

$$\hat{\beta}_S = \underset{\beta}{\operatorname{argmin}} (Y - X\beta)^T (Y - X\beta) + \lambda \|\beta\|_0$$

where $\|\beta\|_0 = \text{number of coefficients} \neq 0$

- Lasso does Shrinkage & Selection.

6) Bayesian Linear Regression

(New material)

← regularization method

We will now consider a Bayesian approach to estimate the parameters of linear models. As mentioned before, a Bayesian approach needs ~~us~~ to specify a prior for the parameters in the model.

Model: $\left. \begin{array}{l} Y | X, \alpha \sim N_N(X\beta, \sigma^2 I) \\ \theta = (\beta, \sigma^2) \sim \text{prior distribution } p(\theta) \end{array} \right\}$

In general, you will use Bayes theorem + numerical method to get a (sample from) the posterior distribution, from which you can get estimators ~~for~~ for β and σ^2 by minimizing the chosen loss function (as seen in one of the numerical examples in the tutorial).

As an example/exercise, we consider here the simple case where:

- σ^2 is known
 - $\beta \sim N_p(0, \sigma_0^2 I_p)$, i.e. zero mean and isotropic
(same variance for all coefficients and no correlations, hence independent)
- ↑ Prior

Remark: The zero-mean prior is not restrictive, you can always transform the variable appropriately. For example, if

$\left(\begin{matrix} \beta_0 \\ \beta_1 \end{matrix} \right) \sim N_2 \left(\begin{pmatrix} \beta \\ \beta \end{pmatrix}, \sigma_0^2 I_2 \right)$, you can transform to:

$$Y_i^* = Y_i - \beta_0 X_i = \underbrace{(\beta_0 - \beta)}_{\beta_0^*} + \underbrace{(\beta_1 - \beta) X_i}_{\beta_1^*} + \epsilon_i$$

$$\Rightarrow \left(\begin{matrix} \beta_0^* \\ \beta_1^* \end{matrix} \right) \sim N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \sigma_0^2 I_2 \right)$$

We have to find posterior distribution. To do this we find joint density of \underline{Y} & β & by conditioning find joint density of \underline{Y} and β : Posterior dist. i.e. $\theta | y$.

$$P(\underline{Y}, \beta | X, \sigma^2, \delta_p^2) = P(\underline{Y} | X, \beta, \sigma^2) P(\beta | \delta_p^2) =$$

$$= \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} \frac{1}{(2\pi\delta_p^2)^{\frac{p}{2}}} \exp \left[-\frac{1}{2\sigma^2} (\underline{Y} - X\beta)^T (\underline{Y} - X\beta) - \frac{1}{2\delta_p^2} \beta^T \beta \right] =$$

$$\propto \exp \left(-\frac{1}{2} \begin{bmatrix} \underline{Y} \\ \beta \end{bmatrix}^T \begin{bmatrix} \frac{1}{\sigma^2} I_N & -\frac{1}{\sigma^2} X \\ -\frac{1}{\sigma^2} X^T & \frac{1}{\delta_p^2} I_p + \frac{1}{\sigma^2} X^T X \end{bmatrix} \begin{bmatrix} \underline{Y} \\ \beta \end{bmatrix} \right)$$

$$\Rightarrow \begin{pmatrix} \underline{Y} \\ \beta \end{pmatrix} \sim N_{N+p} \left(\underline{0}, \Sigma^{-1} \right)$$

From the properties we have seen for the multivariate normal:

$$\beta | \underline{Y} \sim N_p \left(\hat{\beta}, \tilde{\Sigma} \right), \text{ where}$$

$$\hat{\beta} = \Sigma_{22} \Sigma_{11}^{-1} \underline{Y} = - \left(\Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{21} \right) \Sigma_{22} \underline{Y} =$$

$$= - \Sigma_{22}^{-1} \Sigma_{21} \underline{Y} = \left(\frac{1}{\delta_p^2} I_p + \frac{1}{\sigma^2} X^T X \right)^{-1} \frac{1}{\sigma^2} X^T \underline{Y}$$

$$= \left(X^T X + \frac{\sigma^2}{\delta_p^2} I_p \right)^{-1} X^T \underline{Y}$$

AND $\tilde{\Sigma} = \Sigma_{22} - \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{21} = \Sigma_{22}^{-1} = \left(\frac{1}{\delta_p^2} I_p + \frac{1}{\sigma^2} X^T X \right)^{-1}$

Posterior dist.

MAP ESTIMATOR:

$$\hat{\beta} = \mathbb{E}[\beta | Y] = \left(X^T X + \frac{\sigma^2}{\delta_p^2} I_p \right)^{-1} X^T Y$$

The more error in data the more regularization I need to add & more uncertainty of prior the estimator will be closer to prior's estimate. More certain of prior & less of data, the more we are shrinking the parameters.

λ : EQUIVALENT TO THE RIDGE REGRESSION ESTIMATOR

PREDICTIVE DISTRIBUTION:

We now want to predict the response Y^* at a new set of values of the predictors $\underline{x}^* = (\phi_1(\underline{x}^1), \dots, \phi_p(\underline{x}^p))$. In a Bayesian setting, we can compute the PREDICTIVE DISTRIBUTION:

$$P(Y^* | \underline{x}^*, \underline{y}, X, \sigma^2, \delta_p^2) = \int P(Y^* | \underline{x}^*, \beta, \sigma^2) P(\beta | Y, X, \sigma^2, \delta_p^2) d\beta$$

This is an example of Gaussian Linear Model:

Assume $\underline{\beta} \sim N_N(\underline{\mu}, \Lambda^{-1})$ and $U | \underline{\beta} = N_N(A\underline{\beta} + \underline{b}, L)$,

then the marginal density of U is also normal and

$$\mathbb{E}[U] = \int \mathbb{E}[U | \underline{\beta} = \underline{\beta}] p(\underline{\beta}) d\underline{\beta} = \int (A\underline{\beta} + \underline{b}) p(\underline{\beta}) d\underline{\beta} = A\underline{\mu} + \underline{b}$$

$$\text{Cov}(U) = \mathbb{E}[(U - A\underline{\mu} - \underline{b})(U - A\underline{\mu} - \underline{b})^T] =$$

Sub $\rightarrow U = A\underline{\beta} + \underline{b} + \underline{\epsilon}$, $\underline{\epsilon} \sim N(0, L)$ $\perp \underline{\beta}$ independent

$$= \mathbb{E}[(A\underline{\beta} + \underline{b} + \underline{\epsilon} - A\underline{\mu} - \underline{b})(A\underline{\beta} + \underline{b} + \underline{\epsilon} - A\underline{\mu} - \underline{b})^T] =$$

$$= \mathbb{E}[\Lambda(\underline{\beta} - \underline{\mu})(\underline{\beta} - \underline{\mu})^T \Lambda^T] + \mathbb{E}[\underline{\epsilon} \underline{\epsilon}^T] =$$

$$\text{Cov}(U) = A \Lambda^{-1} A^T + L$$

BACK TO THE PREDICTIVE DISTRIBUTION:

$$\left. \begin{aligned} Y^* | \underline{x}^*, \underline{y}, X, \sigma^2, \sigma_p^2 &\sim N_p \left(\underline{\beta}^T \hat{\beta}, \sigma^2 I + \underline{x}^{*T} \tilde{\Sigma} \underline{x}^* \right) \\ \text{Since: } \underline{\beta} | \underline{y}, X, \sigma^2, \sigma_p^2 &\sim N_p \left(\hat{\beta}, \tilde{\Sigma} \right) \\ \text{and } Y^* | \underline{x}^*, \underline{\beta}, \sigma^2 &\sim N \left(\underline{x}^{*T} \underline{\beta}, \sigma^2 \right) \end{aligned} \right\}$$

Uncertainty in $\hat{\beta}$
↓
 $\tilde{\Sigma}$

BAYESIAN MODEL SELECTION:

M_1, \dots, M_L

We now have candidate models Π_1, \dots, Π_L , each with a prior probability $p(\Pi_i)$ of being the "true" model.

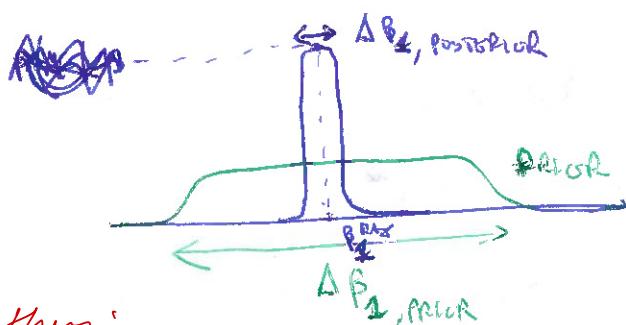
A natural way to do model selection is to look at posterior probability

$$p(\Pi_i | \underline{y}) \propto p(\Pi_i) p(\underline{y} | \Pi_i)$$

MODEL EVIDENCE

$$p(\underline{y} | \Pi_i) = \int p(\underline{y} | \Pi_i, \underline{\beta}_i) p(\underline{\beta} | \Pi_i) d\underline{\beta}$$

Let's try to understand the model evidence in a simple scenario:



In this case then:

$$\text{Then, } p(\underline{y} | \Pi_i) \approx p(\underline{y} | \Pi_i, \underline{\beta}_i^{\text{max}}) \frac{\Delta \beta_i, \text{posterior}}{\Delta \beta_i, \text{prior}}$$

Model evidence Likelihood of observation Max Posterior.

Taking the log:

$$\log P(Y | \Pi_c) \approx \log P(Y | \Pi_c, \beta_{\text{prior}}^{\text{true}}) + \underbrace{\log \left(\frac{\Delta P_{\text{posterior}}}{\Delta P_{\text{prior}}} \right)}_{\text{Posterior ratio}}$$

Log-likelihood in a frequentist setting

If we consider many parameters with the same ratio β_1, \dots, β_p

$$\frac{\Delta \text{Posterior}}{\Delta \text{Prior}},$$

$$\log \left\{ P(Y | \Pi_c) \right\} \approx \log P(Y | \Pi_c, \beta_{\text{prior}}^{\text{true}}) + p \log \left(\frac{\Delta P_{\text{posterior}}}{\Delta P_{\text{prior}}} \right)$$

Model fit

Model fit

↳ Posterior for complexity

However, a better approach would be to do model averaging instead of model selection:

$$\left\{ P(Y^* | \underline{x}, \underline{y}, \underline{x}, \sigma^2, \sigma_p^2) = \sum_{i=1}^L P(Y^* | \underline{x}, \underline{y}, \underline{x}, \sigma^2, \sigma_p^2, \Pi_i) P(\Pi_i | Y) \right\}$$

Remark: In theory, if the true model is one of Π_1, \dots, Π_L , Bayesian approach overfits even when using all the data to fit the model.

- In practice, there are a couple of cases:
 - Prior specification
 - True model not in the set Π_1, \dots, Π_L
- It is better to use a train/test set approach to check the model: forecast predictive checks

End of Week VI

Week VII

- EMPIRICAL BAYES → Get prior from data

Until now, we have assumed σ^2 known, which is unrealistic.

We can of course put a prior on σ^2 as well.

For example : $\frac{1}{\sigma^2} \sim \text{Gamma}(\alpha, \beta)$

$$\beta \sim N_p(0, \frac{\sigma^2}{\alpha} I_p) \text{ independently}$$

(We will see this example in the tutorial)

More complicated priors can be used for (β, σ^2) ,

and numerical methods (NCLC) can be used to sample from the posterior.

However, you still need to choose the hyperparameters $(\alpha, \beta, \sigma_p^2)$

An alternative is to use an approximation, known as

- EMPIRICAL BAYES OR EVIDENCE APPROXIMATION.

Generally speaking, empirical Bayes refers to the fact that the priors (or some of its parameters) are estimated from the data.

Note that this violates the standard Bayes procedure, where the prior needs to be known in advance.

We will again consider a special scenario:

$$\left. \begin{array}{l} \hat{\beta} | \beta, \sigma^2 \sim N_n(X\beta, \sigma^2 I_n) \\ \beta \sim N_p(0, \sigma_p^2 I_p) \\ \sigma^2 \sim p(\sigma^2) \xrightarrow{\text{prior on } \sigma^2} \\ \sigma_p^2 \sim p(\sigma_p^2) \xrightarrow{\text{hyperprior on } \sigma_p^2} \end{array} \right\}$$

Now, the predictive distribution ~~now~~ is:

$$P(\gamma^* | \gamma) = \int \int \int p(\gamma^* | \beta, \sigma^2) P(\beta | \gamma, \sigma_\beta^2, \sigma^2) P(\sigma^2, \sigma_\beta^2 | \gamma) d\beta d\sigma^2 d\sigma_\beta^2$$

since we need to integrate w.r.t. σ^2 and σ_β^2 .

Approximation (Empirical Bayes) 1: The posterior $p(\sigma^2, \sigma_\beta^2 | \gamma)$ is

↳ Assumption ①

SHARPLY PEAKED around $\hat{\sigma}^2, \hat{\sigma}_\beta^2$

↳ $\hat{\sigma}^2, \hat{\sigma}_\beta^2$ values

$$\Rightarrow p(\gamma^* | \gamma) \approx \int p(\gamma^* | \beta, \hat{\sigma}^2) P(\beta | \gamma, \hat{\sigma}^2, \hat{\sigma}_\beta^2) d\beta$$

(Same expression of the predictive distribution ~~now~~ with σ^2 and σ_β^2 known, but plugging in $\hat{\sigma}^2$ and $\hat{\sigma}_\beta^2$ instead)

So what remains is only to find $\hat{\sigma}^2, \hat{\sigma}_\beta^2$.

Approximation (Empirical Bayes) 2: $p(\sigma^2, \sigma_\beta^2) \approx \text{constant}$ ↳ Assumption ②

↳ Can be obtained by
maximising marginal likelihoods below

(Intrinsic prior)

$$\Rightarrow (\hat{\sigma}^2, \hat{\sigma}_\beta^2) = \underset{\sigma^2, \sigma_\beta^2}{\operatorname{argmax}} P(\sigma^2, \sigma_\beta^2 | \gamma) =$$

$$= \underset{\sigma^2, \sigma_\beta^2}{\operatorname{argmax}} P(\gamma | \sigma^2, \sigma_\beta^2) P(\sigma^2, \sigma_\beta^2) = \text{constant}$$

$$= \underset{\sigma^2, \sigma_\beta^2}{\operatorname{argmax}} P(\gamma | \sigma^2, \sigma_\beta^2) = \underset{\sigma^2, \sigma_\beta^2}{\operatorname{argmax}} \int P(\gamma | \beta, \sigma^2) P(\beta | \sigma_\beta^2) d\beta$$

↗ Marginal
likelihood

↗ DENSITY
FUNCTION

→ Likelihood of my data
given only some parameters
as β is not here.

The marginal likelihood $p(\gamma | \sigma^2, \delta_\theta^2)$ can be maximized numerically, either directly, or using an EM algorithm with β as latent variable (see problem sheet 7).

What is interesting to note is this allows us to estimate

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

from the data, without doing any cross-validation.

↳ REGULARISATION PENALTY penalty
in the Bayesian MAP ESTIMATOR

This process is also called Evidence approx. in Machine Learning.

Gaussian Process Regression

This is an extension of Bayesian linear model, where we do not specify the expression of the conditional mean $\underline{H}(\underline{x})$.

Indeed, in a Bayesian linear model we have: $Y_i(x_i) = H(x_i) + \varepsilon_i$, with $H(\underline{x}) = \underline{\beta}^\top \phi(\underline{x})$ and prior $\underline{\beta} \sim N_p(0, \sigma_p^2 I)$.

For N observations,

$$\underline{H} = \begin{bmatrix} H(x_1) \\ \vdots \\ H(x_N) \end{bmatrix} \sim N_N(0, K) \text{ a priori.}$$

because $\underline{H} = X\underline{\beta}$, so it is a Gaussian (linear function of a Gaussian r.v.) with $E[\underline{H}] = X E[\underline{\beta}] = 0$ and $K = \text{Var}(\underline{H}) = X \text{Var}(\underline{\beta}) X^\top = \sigma_p^2 X X^\top$.

Note that $k_{ij} = \sigma_p^2 \phi(x_i) \phi(x_j)^\top$, which is a function of the evaluations of the ~~same~~ predictors in observations i and j .

Therefore, we can redefine the model by using the prior

$$\underline{H} \sim N_N(0, K)$$

This motivates us to generalize this model by changing expectation and variance, i.e. defining a prior directly on the space of the functions $H(\underline{x})$.

A Gaussian Process is a probability distribution over functions $H(x)$, such that the evaluation of H at an arbitrary set of points x_1, \dots, x_N has distribution

$$\begin{bmatrix} H(x_1) \\ \vdots \\ H(x_N) \end{bmatrix} \sim N_N \left(m(x), K^{\text{GP}} \right),$$

↳ GP = Gaussian Process

with $K_{ij} = k(x_i, x_j)$
 ↳ Kernel function

The kernel function $k(\cdot, \cdot)$ needs to be such that K (^{GP}_{matrix}) is positive definite $\forall x_1, \dots, x_N$. & Symmetric.

OBTAINING VALID KERNELS:

Useful rules to construct valid kernels:

- Sum: If $k_1(x, x^*)$ and $k_2(x, x^*)$ are valid kernels, then $k(x, x^*) = k_1(x, x^*) + k_2(x, x^*)$ is a valid kernel.

Proof: $K_{12} = k(x_1, x_2) = k_1(x_1, x_2) + k_2(x_1, x_2)$

$$\Rightarrow K = k_1 + k_2$$

$$v^T K v = \underbrace{v^T k_1 v}_{\geq 0} + \underbrace{v^T k_2 v}_{\geq 0} \geq 0$$

- Product: If $k_1(\cdot, \cdot)$ and $k_2(\cdot, \cdot)$ are valid kernels, then $k(x, x^*) = k_1(x, x^*)k_2(x, x^*)$ is valid.
- Proof: Let $X_1(x)$ be a sample from $\text{GP}(X_1 | 0, k_1)$ Gaussian process
 $X_2(x) \sim \text{GP}(X_2 | 0, k_2)$, independently
 $k_{12} = \mathbb{E}[X_1(x_1)X_1(x_2)X_2(x_3)X_2(x_4)] =$
 $X_1 \perp X_2 = \mathbb{E}[X_1(x_1)X_1(x_2)]\mathbb{E}[X_2(x_3)X_2(x_4)] = (k_1)_{12}(k_2)_{12}$
 Thus is a covariance matrix, so $k(\cdot, \cdot)$ is a valid kernel
 (because K is positive semi-definite)

- Variable Transformation: $k_1(\cdot, \cdot)$ valid,

$$k(x, x^*) = k_1(f(x), f(x^*)) \text{ is valid by using linear transformation } f \text{ in } f.$$

- $k(x, x^*) = g(x)g(x^*)$ is valid
- $k(x, x^*) = \exp(k_1(x, x^*))$ is valid
- $k(x, x^*) = C(k_1(x, x^*))^\lambda$, $C > 0, \lambda > 0$, is valid
- EXAMPLES:

- Gaussian Kernel: $k(x, x^*) = \exp\left(-\frac{\|x - x^*\|^2}{2\sigma^2}\right)$

- Polynomial Kernel: $k(x, x^*) = (x^T x^*)^\lambda$

• GP REGRESSION

Now assume $Y_i = \mu_i + \varepsilon_i$, $\mu_i = \mu(\mathbf{x}_i)$, with

$$\mu(\mathbf{x}) \sim GP(\mu_0, K)$$

$$\varepsilon_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$$

$\mu_0 = 0$
Kernel function
or Covariance
function K

$$\begin{aligned} & \Rightarrow Y_i | \mu \sim N_N(\mu, \sigma^2 I) \\ & \mu \sim N_N(0, K) \leftarrow \text{Prior} \end{aligned}$$

(see last lecture)

This is again an example of GAUSSIAN LINEAR MODEL. Then
the marginal distribution of \mathbf{Y} is Gaussian.

$$\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]] = 0$$

$$\text{Var}(Y) = \cancel{\mathbb{E}[Y^T]} \mathbf{I} K \mathbf{I}^T + \sigma^2 \mathbf{I} = k + \sigma^2 \mathbf{I}$$

$$\Rightarrow \mathbf{Y} \sim N_N(0, k + \sigma^2 \mathbf{I})$$

PREDICTIVE DISTRIBUTION:

$$Y_0 | \mathbf{Y} \sim N\left(K_0^T (k + \sigma^2 \mathbf{I})^{-1} \mathbf{Y}, C_0 - K_0^T (K_0 \sigma^2)^{-1}\right)$$

$$K_0 = \begin{pmatrix} K(x_1, x_0) \\ \vdots \\ K(x_n, x_0) \end{pmatrix}$$

$$C_0 = K_0^T (K + \sigma^2 \mathbf{I})^{-1} K$$

$$C_0 = k(x_0, x_0) + \sigma^2$$

This is obtained by applying the formula of conditional distribution

$$\text{to } \begin{bmatrix} Y_0 \\ \mathbf{Y} \end{bmatrix} \sim N\left(0, C\right) \quad C = \begin{pmatrix} C_0 & K_0^T \\ K_0 & K + \sigma^2 \mathbf{I} \end{pmatrix}$$

Week 8

CLASSIFICATION

(Supervised learning)

The classification problem is a specific type of prediction where the response variable Y is a categorical variable, i.e. it can take value in the set $\{C_1, \dots, C_K\}$.
(Or simply $\{1, \dots, K\}$, but remember they are just labels, there is no order).

What should we do to predict the class at a new observation Y , if we are given data: $(X_{1i}, \dots, X_{pi}, Y_i), i=1, \dots, N$?

DECISION THEORY FOR CLASSIFICATION

We want to choose the classifier (the predictor in this case) by minimizing a loss function. However, the loss functions that we have seen for parameter estimation do not make sense here, since the response variable is categorical.

Let's say $\hat{Y}(x) \in \{C_1, \dots, C_K\}$ is our prediction/classification, then the loss function will have a finite number of values given by all the possible combinations PREDICTED CLASS - TRUE CLASS. Then, it can be written as a $K \times K$ matrix L where

$L_{ij} = \text{cost of predicting } C_j \text{ when the true class is } C_i$

In particular $L_{ii} = 0$ (correct decision), $L_{ij} \geq 0$ if $j \neq i$

Simplest case (0-1 Loss): $L_{ij} = 1 + \delta_{ji}$

→ other diagonal of matrix greater than 0
matrix greater than 0
0 & diagonal = 0

But in practice we may want to have varying costs because some mistakes are worse than others (e.g., healthy patient \rightarrow unhealthy v.s. unhealthy \rightarrow healthy)

Now we minimize the expected loss:

$$\mathbb{E} [L(Y, \hat{Y}(x))] = \mathbb{E}_X [\mathbb{E} [L(Y, \hat{Y}(x)) | X]] = \\ = \mathbb{E}_X \sum_{k=1}^K L(\hat{Y}_k, \hat{Y}(x)) P(Y_k | X)$$

$$\Rightarrow \hat{Y}(x) = \operatorname{arg\,min}_{c_1, \dots, c_K} \sum_{k=1}^K L(\hat{Y}_k, c_k) P(Y_k | X=x)$$

If we assume the 0-1 loss function, this becomes

$$\hat{Y}(x) = \operatorname{arg\,min}_{c_1, \dots, c_K} [1 - P(Y_k = c_k | X=x)]$$

$$\hat{Y}(x) = \tilde{C} \quad \text{where } P(Y=\tilde{C} | X=x) = \max_{c_1, \dots, c_K} P(Y=c | X=x)$$

\Rightarrow This is known as Bayes classifier, i.e. we classify as the most probable class according to the conditional distribution $P(Y | X)$.

(Look up picture on
1-D input from
this course)

In any case, to minimize the expected loss function we need to estimate this conditional distribution.

SYNTHESIS : Approaches to classification

(a) DISCRIMINATIVE APPROACH.

Define a function $f_k(x)$, called a discriminant function, and then. ~~then predict~~. $\hat{Y}(x) = \arg \max_K f_k(x)$

Note: $f_k(x) = P(Y=k | x)$ is a special case, sometimes called DISCRIMINATIVE MODEL

(b) GENERATIVE MODEL

Model the joint distribution $P(X, Y)$ and then classify by conditioning on X .

It is more demanding than approach (a). (Also note that X may have very large dimension)

Its advantage is that we can obtain also $P(x)$, so we will be able to say if a new observation we are trying to classify has low probability under the model

⇒ OUTLIER DETECTION
outlier

8) LOGISTIC REGRESSION

Let us assume we have 2 classes C_1 and C_2 .

$$\Rightarrow Y \sim \text{Bernoulli}(\rho)$$

ρ = probability of belonging to class C_1

$$P(Y=C_1) = \begin{cases} \rho & \text{if } C_1 = C_1 \\ 1-\rho & \text{if } C_1 = C_2 \end{cases}$$

$$\Rightarrow Y|X \sim \text{Bernoulli}(\phi(x)) \quad (\text{THIS IS OUR DISCRIMINATIVE MODEL})$$

$$P(x) = P(C_1|x, \beta) = \frac{1}{1 + \exp(-\phi(x)^T \beta)} \quad (*)$$

↳ LOGISTIC REGRESSION.

This is a special case of GENERALIZED LINEAR MODEL.

There are some good mathematical reason to choose the specific expression in (*), since it is the so-called canonical link (but we are not going to discuss it).

In practice, the only requirement is to be a function from $\mathbb{R} \rightarrow [0, 1]$, so that we can map $\phi(x)^T \beta$ (i.e. a linear combination of our predictors) into something that is a probability.

Other options are:

$$\text{PROBIT REGRESSION : } P(C_1|x, \beta) = \Phi(\phi(x)^T \beta)$$

corresponding

↳ C.D.F. OF A

STANDARDIZED
NORMAL
DISTRIBUTION

LOG-LOG REGRESSION :

$$P(C_1|x, \beta) = \frac{e^{\phi(x)^T \beta}}{1 + e^{\phi(x)^T \beta}}$$

Now come we have the data $(x_{1i}, \dots, x_{di}, y_i) \quad i=1, \dots, n$ and we want to estimate the logistic regression parameters β .

The likelihood is $P(Y|P)$

Coding: $y=1 \text{ if } c_1$
 $y=0 \text{ if } c_2$

$$P(Y|P) = \prod_{i=1}^n P_i^{y_i} (1-P_i)^{1-y_i}, \quad \begin{cases} P(y_i=1) = P_i \\ P(y_i=0) = 1-P_i \end{cases}$$

where $P_i = \frac{1}{1 + \exp(-\phi(x_i)^T P)} = h(\phi(x_i)^T P)$

The log-likelihood is

$$\ell = \sum_{i=1}^n \left[y_i \log P_i + (1-y_i) \log (1-P_i) \right]$$

Computing the gradient: (Gradient of likelihood $\frac{\partial \ell}{\partial \beta_j}$)

$\star \quad \left[\frac{\partial \ell}{\partial \beta_j} = \sum_{i=1}^n y_i \frac{1}{P_i} \frac{\partial P_i}{\partial \beta_j} - (1-y_i) \frac{1}{1-P_i} \frac{\partial P_i}{\partial \beta_j} \right]$

Now, compute term $\frac{\partial P_i}{\partial \beta_j}$:

Now, $\frac{\partial P_i}{\partial \beta_j} = h'(\phi(x_i)^T P) \frac{\partial}{\partial \beta_j} (\phi(x_i)^T P) = h'(\phi(x_i)^T P) \underbrace{\phi_j(x_i)}_{\phi_j(x_i)} =$

$$\begin{aligned} h'(a) &= + \frac{\exp(a)}{(1+\exp(a))^2} = \frac{1}{1+\exp(-a)} \left(1 - \frac{1}{1+\exp(-a)} \right) = \\ &= h(a) (1-h(a)) \end{aligned}$$

$$\Rightarrow \left\{ \frac{\partial P_i}{\partial \beta_j} = h(\phi(x_i)^T P) [1 - h(\phi(x_i)^T P)] \underbrace{\phi_j(x_i)}_{\phi_j(x_i)} = \right.$$

$$= P_i (1-P_i) \phi_j(x_i) \quad \rightarrow \text{sub this in } \star$$

$$\Rightarrow \frac{\partial \ell}{\partial \beta_j} = \sum_{i=1}^n \left[\frac{y_i}{P_i} P_i (1-P_i) \phi_j(x_i) - (1-y_i) \frac{1}{1-P_i} P_i (1-P_i) \phi_j(x_i) \right] =$$

$$= \sum_{i=1}^n (y_i - P_i) \overbrace{\phi_j(x_i)}^{x_{ij}} = T(Y-P^T X_j) \quad P = T_{Pj}^T$$

gradient of likelihood

I can write this as vector incarnation of my gradient as:

$$\Rightarrow \left[\nabla_{\beta} \ell = X^T (\underline{y} - \underline{\phi}) \right] \quad \text{Note: For a linear model, we get } \nabla \ell = X^T (\underline{y} - X \underline{\beta})$$

The Hessian of ℓ is then: (Non-linear β & since we do have nice linear expression of ℓ in terms of β)
(indefinite w.r.t $\underline{\beta}$)

$$\begin{aligned} \frac{\partial^2 \ell(\beta)}{\partial \beta_k \partial \beta_j} &= \frac{\partial}{\partial \beta_k} \frac{\partial \ell}{\partial \beta_j} = \frac{\partial}{\partial \beta_k} \left[\sum_{i=1}^n \phi_j(x_i) (y_i - \beta_i) \right] = \\ &= \sum_{i=1}^n \phi_j(x_i) \frac{\partial \beta_i}{\partial \beta_k} = \sum_{i=1}^n \phi_j(x_i) p_i (1-p_i) \phi_k(x_i) = \\ &= \sum_{i=1}^n -x_{ik} p_i (1-p_i) x_{ij} = (-X^T R X)_{kj} \\ R &= \begin{pmatrix} p_1(1-p_1) & & \\ & \ddots & \\ & & p_n(1-p_n) \end{pmatrix} \end{aligned}$$

The Hessian is negative-definite because, for $\underline{v} \neq 0$,

$$\underbrace{\underline{v}^T X^T R X \underline{v}}_{\underline{u}^T \underline{u}} = \underline{u}^T R \underline{u} = \sum_{i=1}^n u_i^2 p_i (1-p_i) > 0$$

$$\Rightarrow \cancel{\text{Hessian}} \quad \text{H}(\ell(\beta)) < 0$$

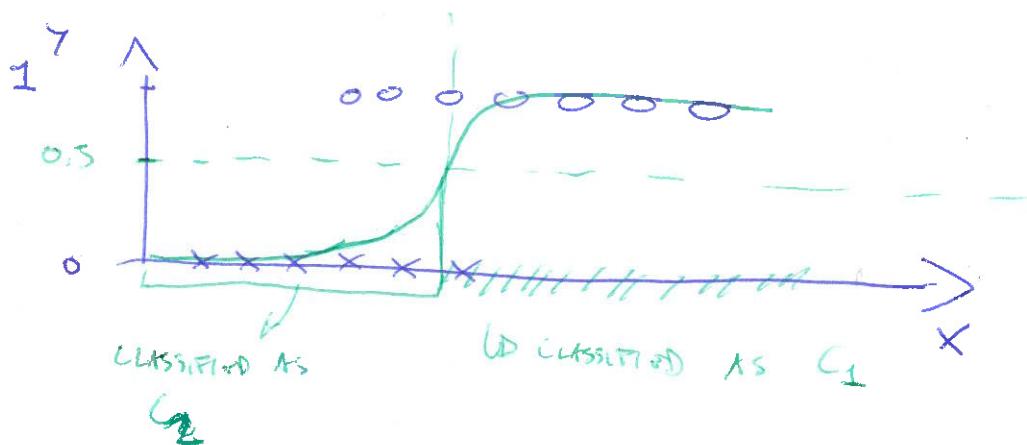
Checking Hessian is negative-definite so, has our maximum.

$\ell(\beta)$ has a unique maximum, so we can use Newton-Raphson method to numerically solve

$$\boxed{\nabla_{\beta} \ell = 0}$$

In this context, it is called Iterative RE-WEIGHTED LEAST SQUARES algorithm because each step is equivalent to solve a weighted least squares problem.

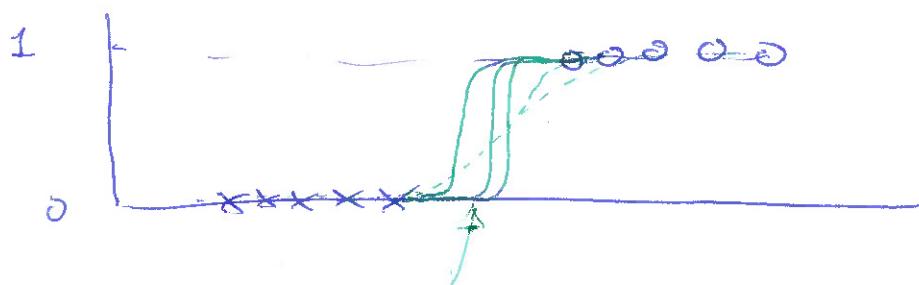
Final representation: For 1D input x



the 0.5 cut-off or what is reported by the Bayes classifier if the loss function is 0-1.

Remark: • There are identifiability issues when data are perfectly separated.

Remove
for page
before



Any of THESE HAVE THE SAME LIKELIHOOD

- Logistic regression can be generalized to K classes using MULTINOMIAL LOGISTIC REGRESSION
- observation $Y \sim \text{Multinomial}(1, x_1, \dots, x_k) \pi$ (where π_i sum to 1)
- $$x_j | X = P(Y = C_j | X) = \frac{e^{\phi(x)^T \beta_j}}{1 + \sum_{k=1}^{K-1} e^{\phi(x)^T \beta_k}}, j=1, K-1$$
- $$\pi_k | X = \frac{1}{1 + \sum_{j=1}^{K-1} e^{\phi(x)^T \beta_j}} \quad (\text{so that } \sum_{j=1}^K \pi_j = 1)$$

BAYESIAN LOGISTIC REGRESSION :

~~Prior~~ Prior : $\beta \sim N(0, \sigma_p^2 I)$

Posterior $\rightarrow p(\beta | y) \propto p(y | \beta) p(\beta)$

by $p(y | \beta) = \ell(\beta) - \frac{1}{2\sigma_p^2} \beta^T \beta + \text{const.}$

$$\ell(\beta) = \sum_{i=1}^n (y_i \ln p_i + (1-y_i) \ln (1-p_i)) - \frac{1}{2\sigma_p^2} \beta^T \beta$$

this is not a Gaussian distribution.

However, it can be shown to be unimodal, so it is common ^{practice} to approximate it with a Gaussian density with the same ^{Mode}. (LAPLACE approximation)

Alternatively, you can use numerical methods.

1-D LAPLACE APPROXIMATION

~~Prior~~ Prior : Consider a generic density

$$p(\beta) = \frac{\exp(\hat{\ell}(\beta))}{Z}$$

Taylor expansion around β_0 : \rightarrow UNIQUES MAXIMUM

$$\beta_0 = \max_{\beta} \hat{\ell}(\beta)$$

$$\left\{ \hat{\ell}(\beta) \approx \hat{\ell}(\beta_0) + \frac{1}{2} \frac{\partial^2}{\partial \beta^2} [\hat{\ell}(\beta)]_{\beta=\beta_0} (\beta - \beta_0)^2 \right\} \left(\begin{array}{l} \text{first term in Taylor expansion} \\ \text{Normalized constant} \\ \text{because } \frac{\partial}{\partial \beta} \hat{\ell}(\beta) = 0 \end{array} \right)$$

$$\Rightarrow p(\beta) \approx \frac{e^{\hat{\ell}(\beta)}}{Z} \exp\left(+\frac{1}{2} \frac{\partial^2}{\partial \beta^2} [\hat{\ell}(\beta)]_{\beta=\beta_0} (\beta - \beta_0)^2\right)$$

Exponential of $\hat{\ell}(\beta)$

$$\Rightarrow \beta \sim N\left(\beta_0, -\frac{1}{\frac{\partial^2}{\partial \beta^2} [\hat{\ell}(\beta)]_{\beta=\beta_0}}\right)$$

P-Dimension Laplace Approximation:

$$\beta \sim N(\beta_0, A^{-1}),$$

where $A = -\Delta \hat{L}(\beta) \Big|_{\beta=\beta_0}$

We have $\hat{L}(\beta) = L(\beta) + \frac{1}{2\sigma^2} \beta^\top \beta$, therefore

$$-\Delta \hat{L}(\beta) = X^\top R X + \frac{1}{\sigma^2} I_p$$

$$R = \begin{pmatrix} \pi_1(1-\pi_1) & & 0 \\ \cdots & \cdots & \cdots \\ 0 & \cdots & \pi_N(1-\pi_N) \end{pmatrix}$$

$$\Rightarrow \left[\beta | y \sim N(\beta_{MAP}, S_N) \right] \quad \textcircled{*}$$

where $S_N = -\Delta \hat{L}(\beta) \Big|_{\beta=\beta_{MAP}}$

Therefore, we can approximate the predictive distribution:

$$\begin{aligned} p(Y=c_1 | x_0, y) &= \int p(Y=c_1 | x_0, \beta) p(\beta | y) d\beta \\ &\simeq \int h(\phi(x_0)^\top \beta) p^\otimes(\beta | y) d\beta \end{aligned}$$

This needs to be computed numerically.

End of week 8

Week 9

10.) Generative Models for Classification

Marginal
 Probability
 → DENSITY
 Function
 P

Strategy: 1) Model the joint probability $P(X, Y) = P(X|Y)P(Y)$

2) Compute the conditional prob. $P(Y|X) = \frac{P(X|Y)P(Y)}{\sum_k P(X|Y=c_k)P(Y=c_k)}$

3) [Assuming a 1-0 loss function] Bayes Classification

$$\hat{Y}(x) = C \quad \text{where} \quad C = \operatorname{argmax}_{c_1, \dots, c_K} P(Y=c_k | X=x)$$

(Alternatively, minimize the expected loss)

$$\hat{Y}(x) = \operatorname{argmin}_{c \in \{c_1, \dots, c_K\}} \sum_{k=1}^K L(c_k, c) P(Y=c_k | X=x)$$

•) Optionally, check $P(X=x)$ is not too small
(Outlier Detection)

Example:

2 classes ; $\{C_1, C_2\}$

Predictor given it belongs to one of 2 classes

$$X|Y=c_k \sim N_p(\mu_k, \Sigma_k), \quad k=1, 2$$

$$P(Y=c_k) = \pi_k \quad (\pi_1=\pi, \pi_2=1-\pi)$$

To compare the two classes, we can simply look at the logarithm of the Ratio between $P(Y=c_1|x)$ and $P(Y=c_2|x)$:

$$\log \frac{P(Y=c_1|x)}{P(Y=c_2|x)} \begin{cases} > 0 & \hat{Y}(x) = C_1 \rightarrow \text{Predicting } C_1 \text{ if } > 0 \\ \leq 0 & \hat{Y}(x) = C_2 \rightarrow \text{.. } C_2 \text{ if } \leq 0 \end{cases}$$

Now,

$$\log \frac{P(Y=c_2|x)}{P(Y=c_1|x)} = \log \frac{P(X|Y=c_2)}{P(X|Y=c_1)} + \log \frac{P(Y=c_2)}{P(Y=c_1)} =$$

$$\log \frac{f_2(x)}{f_1(x)} + \log \frac{\pi_2}{\pi_1} =$$

$$\log \frac{\pi_2}{\pi_1} = \log \frac{\pi}{1-\pi}$$

$$= \log f_1(x) - \log f_2(x) + \log \frac{\pi_1}{\pi_2} =$$

where $f_k(\cdot)$ is the p.d.f. of $N_p(\underline{\mu}_k, \Sigma_k)$

$$= -\frac{1}{2} \log |\Sigma_1| + \frac{1}{2} \log |\Sigma_2| + \log \frac{\pi_1}{\pi_2} +$$

$$- \frac{1}{2} (\underline{x} - \underline{\mu}_1)^T \Sigma_1^{-1} (\underline{x} - \underline{\mu}_1) +$$

$$+ \frac{1}{2} (\underline{x} - \underline{\mu}_2)^T \Sigma_2^{-1} (\underline{x} - \underline{\mu}_2) =$$

$$= -\frac{1}{2} \underbrace{\log |\Sigma_1| + \frac{1}{2} \log |\Sigma_2| + \log \frac{\pi_1}{\pi_2}}_{\beta_0} - \frac{1}{2} \underline{\mu}_1^T \Sigma_1^{-1} \underline{\mu}_1 + \frac{1}{2} \underline{\mu}_2^T \Sigma_2^{-1} \underline{\mu}_2 +$$

$$+ \underbrace{\underline{\mu}_1^T \Sigma_1^{-1} \underline{x} - \underline{\mu}_2^T \Sigma_2^{-1} \underline{x}}_{\beta^T \underline{x}}$$

\uparrow
terms without \underline{x}

\leftarrow terms in \underline{x}

\leftarrow term quadratic in \underline{x}

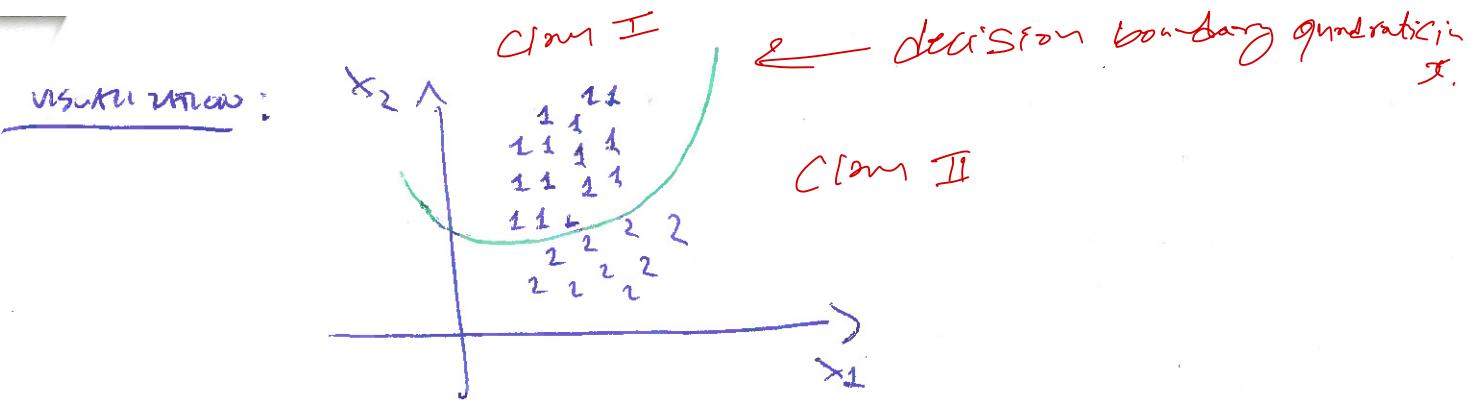
$$\star -\frac{1}{2} \underline{x}^T \Sigma_1^{-1} \underline{x} + \frac{1}{2} \underline{x}^T \Sigma_2^{-1} \underline{x} =$$

$$= \boxed{\beta_0 + \beta^T \underline{x} + \frac{1}{2} \underline{x}^T (\Sigma_2^{-1} - \Sigma_1^{-1}) \underline{x}}$$

QUADRATIC DISCRIMINANT FUNCTION

This method is called QUADRATIC DISCRIMINANT ANALYSIS (QDA)

Remark: Starting with a Generative approach, we ended up with a discriminant function we could have chosen without any probabilistic model



Special case: LINEAR DISCRIMINANT ANALYSIS (LDA) \hookrightarrow Special Case of QDA

If $\Sigma_2 = \Sigma_1 = \Sigma$:

$$\text{Lg} \frac{P(Y=c_1|x)}{P(Y=c_2|x)} = \text{Lg} \frac{\pi_1}{\pi_2} - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)^T \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) +$$

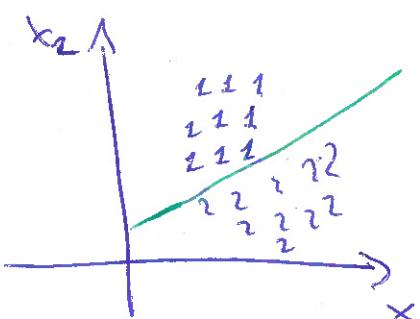
$\underbrace{\phantom{\text{Lg} \frac{\pi_1}{\pi_2}}}_{\beta_0}$

$$+ (\underline{\mu}_1 - \underline{\mu}_2)^T \Sigma^{-1} \underline{x} =$$

$\underbrace{\phantom{(\underline{\mu}_1 - \underline{\mu}_2)^T \Sigma^{-1} \underline{x}}}_{\beta^T \underline{x}}$

$$= \boxed{\beta_0 + \beta^T \underline{x}}$$

LINEAR
DISCRIMINANT
FUNCTION



Remark: • LDA & QDA are often used also when the data are not Gaussian, simply because linear and quadratic boundaries work well as discriminative functions. However, it may make more sense to ~~use~~ use different ~~additive~~ (non-additive) expressions for β_0 and β_1 . • Let w in LDA, $w^T x \rightarrow$ quadratic boundary predictor

Parameter estimation for LDA

(QPT left for exercise)

$$(\underline{x}_l, y_l), l=1, \dots, N$$

(Just like Maximum Likelihood estimation)

$$Y_l = \begin{cases} 1 & \text{if } \text{l-th observation belongs to } C_1 \\ 0 & \text{else} \end{cases}$$

LDA model is equivalent to a Gaussian mixture with 2 components where complete data likelihood is available:

$$P(y, x | \pi, \mu_1, \mu_2, \Sigma) = \prod_{l=1}^N \left[\pi f_1(\underline{x}_l) \right]^{y_l} \left[(1-\pi) f_2(\underline{x}_l) \right]^{1-y_l}$$

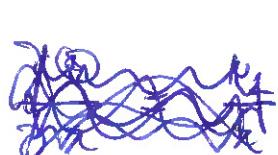
$$\log P(y, x | \pi) \propto \sum_{l=1}^N y_l \left(\log \pi + \frac{1}{2} \log |\Sigma| + -\frac{1}{2} (\underline{x}_l - \underline{\mu}_1)^T \Sigma^{-1} (\underline{x}_l - \underline{\mu}_1) \right) +$$

$$+ (1-y_l) \left(\log (1-\pi) - \frac{1}{2} \log |\Sigma| + -\frac{1}{2} (\underline{x}_l - \underline{\mu}_2)^T \Sigma^{-1} (\underline{x}_l - \underline{\mu}_2) \right)$$

max. w.r.t π first

$$\therefore \max_{\pi} \log P(y, x | \pi) = \max_{\pi} \underbrace{\sum_{l=1}^N y_l \log \pi + (1-y_l) \log (1-\pi)}_{\textcircled{*}}$$

$$N_1 = \sum_{l=1}^N y_l \quad \rightarrow \max_{\pi} \underbrace{N_1 \log \pi + (N-N_1) \log (1-\pi)}$$



$$\frac{\partial \textcircled{*}}{\partial \pi} = \frac{N_1}{\pi} - \frac{(N-N_1)}{1-\pi} = 0$$

$$\Rightarrow \frac{\pi}{1-\pi} = \frac{N_1}{N-N_1} \Rightarrow \boxed{\pi = \frac{N_1}{N}}$$

taking derivative &
setting equal to zero
and in MLE.

Now maximize wrt μ_1 :

$$\bullet) \max_{\mu_1} \log(\rho(y_i, x_i | \theta)) = \max_{\mu_1} \sum_{i=1}^n y_i \left[-\frac{1}{2} (\underline{x}_i - \underline{\mu}_1)^T \Sigma^{-1} (\underline{x}_i - \underline{\mu}_1) \right] =$$

$$= \max_{\mu_1} \sum_{j=1}^{N_1} -\frac{1}{2} (\underline{x}_j^{(1)} - \underline{\mu}_1)^T \Sigma^{-1} (\underline{x}_j^{(1)} - \underline{\mu}_1)$$

↙ Sum over
only one for
which $y_j = 1$

where $\underline{x}_j^{(1)}$ vs the j -th obs for which $y_j = 1$

$$\Rightarrow \text{From Derivative Normal Results } : \left[\hat{\mu}_1 = \frac{1}{N_1} \sum_{j=1}^{N_1} \underline{x}_j^{(1)} \right] \quad \text{↙ We have done this for Gaussian Mixture}$$

$$\text{Analogously: } \left[\hat{\mu}_2 = \frac{1}{N_2} \sum_{j=1}^{N_2} \underline{x}_j^{(2)} = \frac{1}{N - N_1} \sum_{j=1}^{N - N_1} \underline{x}_j^{(2)} \right]$$

Maximize wrt. Σ :

$$\bullet) \max_{\Sigma} \log(\rho(y_i, x_i | \theta)) =$$

$$= \max_{\Sigma} \sum_{i=1}^n \left[-\frac{1}{2} \log |\Sigma| - \frac{y_i}{2} (\underline{x}_i - \hat{\mu}_1)^T \Sigma^{-1} (\underline{x}_i - \hat{\mu}_1) + \right. \\ \left. + -\frac{1-y_i}{2} (\underline{x}_i - \hat{\mu}_2)^T \Sigma^{-1} (\underline{x}_i - \hat{\mu}_2) \right] =$$

$$= \max_{\Sigma} -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \text{trace} \left(\Sigma^{-1} \sum_{j=1}^{N_1} (\underline{x}_j^{(1)} - \hat{\mu}_1) (\underline{x}_j^{(1)} - \hat{\mu}_1)^T \right) +$$

$$- \frac{1}{2} \text{trace} \left(\Sigma^{-1} \sum_{j=1}^{N - N_1} (\underline{x}_j^{(2)} - \hat{\mu}_2) (\underline{x}_j^{(2)} - \hat{\mu}_2)^T \right) =$$

$$= \max_{\Sigma} -\frac{N}{2} \log |\Sigma| - \frac{1}{2} \text{trace} (\Sigma^{-1} V)$$

$$\text{where } V = \frac{1}{N} \left[\sum_{j=1}^{N_1} (\underline{x}_j^{(1)} - \hat{\mu}_1)^T (\underline{x}_j^{(1)} - \hat{\mu}_1) + \sum_{j=1}^{N - N_1} (\underline{x}_j^{(2)} - \hat{\mu}_2)^T (\underline{x}_j^{(2)} - \hat{\mu}_2) \right]$$

from multivariate Normal results

Then estimates follow

$$\hat{\Sigma} = V$$

Remark: THE UNBIASED ESTIMATOR

$$\hat{\Sigma} = \frac{1}{N-2} \left(\sum_{j=1}^{N_1} (x_j^{(1)} - \hat{x}_1)(x_j^{(1)} - \hat{x}_1)^T + \sum_{j=1}^{N-N_1} (x_j^{(2)} - \hat{x}_2)(x_j^{(2)} - \hat{x}_2)^T \right)$$

is sometimes preferred,
Sometimes preferred

- Further remark on this written after
Page (Multiclass LDA/QDA).

MULTI-CLASS LDA / QDA

$$x | Y = C_k \sim N(\mu_k, \Sigma_k), k=1, \dots, K$$

We describe the decision rules we obtain for any 2 pairs it does as a discriminative function; $\max_{\max} \log P(Y=C_k | x)$:

QDA: $f_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$

LDA: $f_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$

then we do classify by maximizing the discriminative functions:

$$\hat{Y}(x) = \arg \max_{k=1, \dots, K} f_k(x)$$

$D_K(x)$

↓
Discriminative
function

• Remark

Remark

- Both LDA and QDA require Σ^{-1} .

$P = \text{no. of predictor}$
 $N = \text{no. of observation}$

Therefore we need $\hat{\Sigma}$ to be invertible, i.e.

$$\text{rank}(\hat{\Sigma}) = p \Rightarrow N > p$$

Therefore LDA and QDA do not work in HIGH-DIMENSIONAL SETTINGS where $p > N$ and, unless some dimensional reduction is performed first.

To estimate parameters,

- You need a genuine sample from (X, Y) , it does not work if we select X values as in linear regression.

joint dist. of $X \times Y$

not work

• OTHER APPROACHES TO GENERATIVE CLASSIFICATION : (Non-Exhaustive Below)

1) REGULARIZED Discriminant Analysis

The covariance of the 2 groups is assumed to be different, but the estimates are regularised so that they are not too different:

$$\hat{\Sigma}_k(x) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}$$

↓
QDA ESTIMATES

LDA ESTIMATE

In between
LDA & DA.

α is chosen by cross-validation.

2) Kernel Density Classification

The density $f_k(x) = P(X | Y = C_k)$ is estimated from the data using Kernel Density Estimation:

$$\hat{f}_k(x) = \frac{1}{N_k} \sum_{j=1}^{N_k} \frac{k_\lambda(x, x_j^{(k)})}{\lambda^p}$$

where $k_\lambda(\cdot, \cdot)$ is a kernel function, for example ~~the Gaussian kernel~~.

The Gaussian kernel: $k_\lambda(x, y) = \frac{1}{(2\pi)^{p/2}} \exp\left(-\frac{1}{2} \frac{\|x-y\|^2}{\lambda^2}\right)$

then,

$$P(Y = C_k \mid X = x_0) = \frac{\hat{\pi}_k \hat{f}_k(x_0)}{\sum_{j=1}^J \hat{\pi}_j \hat{f}_j(x_0)}$$

Naive Bayes Classification

Since ~~the~~ density estimator does not work well in high dimension,
(see the curse of dimensionality)

a popular alternative is to approximate the multivariate density
with:

$$f_k(x) = \prod_{j=1}^p f_{kj}(x_j)$$

(i.e., assuming independence of the predictors)

this works well because it dramatically decreases the variance of
the estimators, even if we are including some bias due to the
incorrect approximation
incorrect

End of Week 9

Week 10

1) Separating Hyperplanes and SVM

We have seen in the past lectures different methods to estimate linear boundaries between decision regions, based on probabilistic models either for (x, y) or $y|x$.

In alternative approach we to ~~estimate~~ learn the "best" linear decision boundaries from the data, without assuming a model for our observations. For simplicity, we ~~restrict~~ our discussion to the 2-class case.

Separating hyperplanes

A separating hyperplane is a linear decision boundary

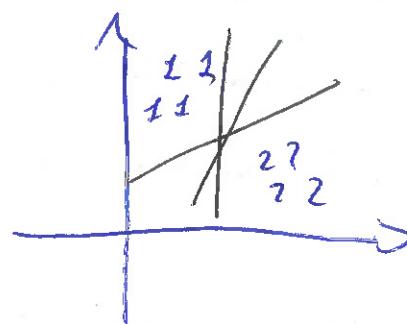
$$\underline{\phi(x)^T \beta}, \text{ such that } \underline{\phi(x_i)^T \beta} = \begin{cases} > 0 & \text{if } y_i = c_1 \\ < 0 & \text{if } y_i = c_2 \end{cases}$$

The number is then ~~$\phi(x_0)^T \beta$~~ $\gamma(x_0) = \text{sign}(\phi(x_0)^T \beta)$.

Note that: a). Data may not be separable for the chosen feature (predictor) space.

problem }

b). If data are separable, the separating hyperplane is not unique



Solution to a) may be choosing a different set of predictors, or using SVM.

Solution to b) is selecting the optimal separating hyperplane.

Optimal separating hyperplane

Critition we want to maximise

$$\hat{f} = \text{argmax } M \leftarrow \begin{array}{l} \text{some number} \\ \text{called Margin} \end{array}$$

$$\|\beta\| = 1$$

$$y_i \cdot \phi(x_i)^T \beta \geq M, \quad i=1, \dots, N$$

$$y_i = \begin{cases} 1 & \text{if } c_1 \\ -1 & \text{if } c_2 \end{cases}$$

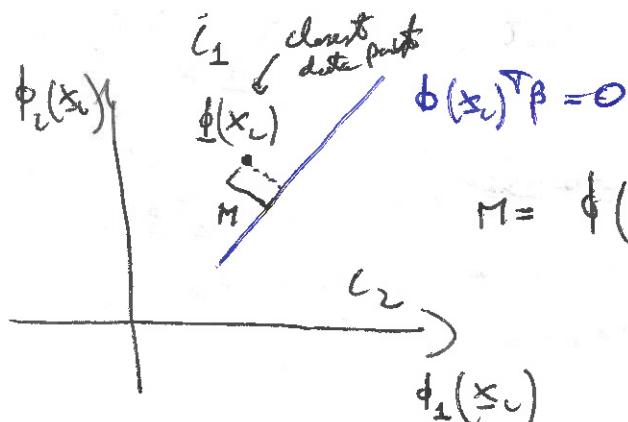
Coding

This set of conditions imply that all points are at distance at least M from the decision boundary.

M is called the MARGIN.

($\|\beta\|=1$ is required because the hyperplane does not change if we rescale the coefficients, i.e. $\phi(x_i)^T \beta = 0 \Leftrightarrow \phi(x_i)^T \lambda \beta = 0$ ($\lambda \neq 0$))

2D example



$M = \phi(x_i)^T \beta$ is the smallest of this is largest possible one for each data point

- $y_i \cdot \phi(x_i)^T \beta \geq M$ guarantees that the distance of the closest point from the boundary (on either side) is largest or equal to M ($\phi(x_i)^T \beta \leq M$ if it lies closer from c_2)

- We can get rid of the $\|\beta\|=1$ constraint by rescaling:

$$\frac{1}{\|\beta\|} (\phi(x_i)^T \beta) \geq M \quad \text{or} \quad \phi(x_i)^T \beta \geq M \|\beta\|$$

and then arbitrarily choose the norm of β , $\|\beta\| = \frac{1}{M}$

We get the equivalent problem:

$$\left. \begin{array}{l} \hat{\beta} = \underset{\beta}{\operatorname{arg\min}} \quad \frac{1}{2} \|\beta\|^2 \\ \gamma_i [\phi(\underline{x}_i)^T \beta] \geq 1, \quad i=1, \dots, N \end{array} \right\} \textcircled{*} \quad \text{becomes minimization problem}$$

Solving problem $\textcircled{*}$:

It is a constrained optimization problem, so we use Lagrange multipliers:

$$\left(\text{primary problem} \right) L = \frac{1}{2} \|\beta\|^2 - \sum_{i=1}^n \lambda_i [\gamma_i \phi(\underline{x}_i)^T \beta - 1]$$

+ KKT conditions

↗ $\lambda_i \geq 0$
 ↗ λ_i range
 ↗ λ_i range multiplication

$$\frac{\partial L}{\partial \beta_j} = 0 \Rightarrow \boxed{\beta_j = \sum_{i=1}^n \lambda_i \gamma_i \phi_j(\underline{x}_i)}$$

~~for β~~

this leads to the Wolfe dual problem

$$(\max) \quad L_d = \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \lambda_i \lambda_k \gamma_i \gamma_k \phi(\underline{x}_i)^T \phi(\underline{x}_k)$$

$$\text{subject to } \lambda_i \geq 0, \quad \sum_{i=1}^n \lambda_i \gamma_i = 0, \quad i=1, \dots, n$$

↳ easier to solve numerically

Lagrangian dual

LKT conditions :

(Lagrange-Kuhn-Tucker)

Gradient of Lagrangian

$$\nabla L = 0 \Rightarrow \beta - \sum_{i=1}^n \lambda_i y_i \phi(x_i) = 0$$

$$y_i \phi(x_i)^T \beta - 1 \geq 0 \quad \text{constraint}$$

$$\lambda_i [y_i \phi(x_i)^T \beta - 1] \geq 0 \quad (\#)$$

$(\lambda_i > 0)$ ↗
Lagrange multiplier

Condition (#) implies that either $\lambda_i = 0$ or

$$y_i \phi(x_i)^T \beta = 1$$

↳ this means the i -th obs
is one of the closest
ones to the boundary,
since its distance is equal
to the MARGIN

Therefore:

$$y_i \phi(x_i)^T \beta = 1$$

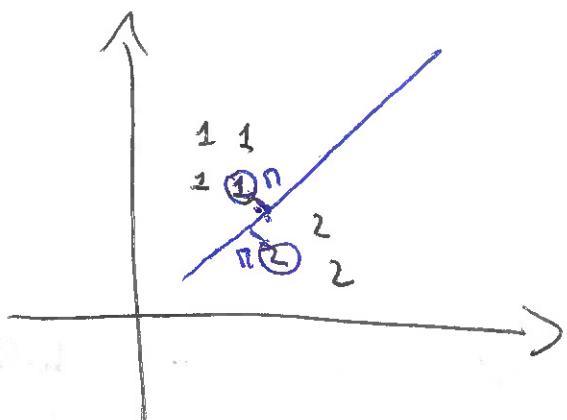
If $\lambda_i > 0$, the i -th obs. is on the margin.
(it lies on the margin)
is close to the boundary and it impacts the expression of the
optimal hyperplane $\rightarrow x_i$ is called support point



If $y_i (\phi(x_i)^T \beta) > 1$, the observation is further away from the
boundary, $\lambda_i = 0 \Rightarrow$ this observation does not impact the
expression of the boundary (since $\beta = \sum_{i=1}^n \lambda_i y_i \phi(x_i)$)

Point: • Since the dual problem depends on $\phi(x_1)^\top \phi(x_2) = k(x_1, x_2)$, and only a small fraction of the data points determines the solution, this approach is referred to as a sparse kernel method.

2D example:



- The decision is $\hat{Y}(x) = \text{sign}(\phi(x)^\top \hat{\beta})$

SUPPORT VECTOR METHODS (for classification)

The optimal hyperplane is a special case of support vector classification.

It is so called because the classification rule depends on a few support vectors (i.e. the ones lying on the margin).

The general formulation allows for the classes to overlap (a little). So now we want to maximize M , but allowing some points to be on the "wrong" side of the boundary.

We therefore reformulate the optimization problem in the following way:

$$\max_{\beta} M$$

$\|\beta\| = 1$

$$y_i (\phi(x_i)^T \beta) \geq M (1 - \varepsilon_i) + v$$

$$\varepsilon_i \geq 0$$

$$\sum_{i=1}^n \varepsilon_i \leq C \quad \leftarrow \text{additional parameter}$$

this means we are allowing data points to be a little bit closer than M to the margin, as long as the sum of the slack variables ε_i is not too large.

As before, we can rewrite this problem as:

$$\hat{\beta} = \arg \min_{\beta} \frac{1}{2} \|\beta\|^2$$

$$y_i (\phi(x_i)^T \beta) \geq 1 - \varepsilon_i$$

$$\varepsilon_i \geq 0$$

$$\sum_{i=1}^n \varepsilon_i \leq C$$

Solving the SVM Problem :

We can write the Lagrangian:

$$L_p = \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n \varepsilon_i - \sum_{i=1}^n \lambda_i [\gamma_i (\phi(x_i)^T \beta) - (1 - \varepsilon_i)] - \sum_{i=1}^n \lambda_i \varepsilon_i$$

$\underbrace{\sum_{i=1}^n \varepsilon_i \leq C}$ $\underbrace{\gamma_i (\phi(x_i)^T \beta) \geq (1 - \varepsilon_i)}$ $\underbrace{\varepsilon_i \geq 0}$

Lagrangian conditions :

- $\nabla L_p = 0 \Rightarrow \hat{\beta} = \sum_{i=1}^n \lambda_i \gamma_i \phi(x_i)$
- $\lambda_i = -\lambda_i \quad \lambda_i [\gamma_i \phi(x_i)^T \beta - (1 - \varepsilon_i)] = 0$
- $\lambda_i \varepsilon_i = 0 \quad \gamma_i \phi(x_i)^T \beta - (1 - \varepsilon_i) \geq 0$

We get that $\lambda_i > 0 \Leftrightarrow \gamma_i \phi(x_i)^T \beta = 1 - \varepsilon_i$

If $\hat{\varepsilon}_i = 0 \Rightarrow$ and $\gamma_i \phi(x_i)^T \beta = 1 \Rightarrow$ ~~$0 < \lambda_i < C$~~

If $\hat{\varepsilon}_i > 0$ and $\gamma_i \phi(x_i)^T \beta = 1 - \hat{\varepsilon}_i \Rightarrow$ ~~$\lambda_i = 0$~~ $\lambda_i = C$

The problem needs to be solved numerically (usually through the simpler dual problem)

Remarks: • SVM can be generalized with different kernels:

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2) \quad \text{Gaussian}$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}^T \mathbf{x})^d \quad \text{Polynomial}$$

and the classifier rewritten as

$$Y(\mathbf{x}) = \underbrace{\phi^T(\mathbf{x}) \beta}_{\text{Sign}} = \underbrace{\phi^T(\mathbf{x})}_{\text{Sign}} \left(\sum_{i=1}^n d_i y_i \phi(\mathbf{x}_i) \right) =$$

$$= \text{Sign} \left(\sum_{i=1}^n d_i y_i \phi(\mathbf{x}_i)^T \phi(\mathbf{x}) \right) =$$

$$\hookrightarrow \text{Sign} \left(\sum_{i=1}^n d_i y_i k(\mathbf{x}, \mathbf{x}_i) \right)$$

- The constant C is playing the role of a regularization parameter and needs to be chosen (for example by cross-validation)
- The cost constant C depends on how large your model is (either in terms of feature space $\phi_i(\mathbf{x}) - \phi_j(\mathbf{x})$ or complexity of the kernel)
- Some people make a distinction between soft SVM and hard SVM (i.e. separating hyperplanes)

End of Week 10