

# Gaussian Models

Congyuan Duan

School of Mathematics  
Sun Yat-sen University

- Introduction
- Why Gaussian Assumption?
- MLE for Gaussian

# 1. Introduction

If  $X \sim N(\boldsymbol{\mu}, \Sigma)$ , the probability density function (pdf) of  $X$  is defined as follow:

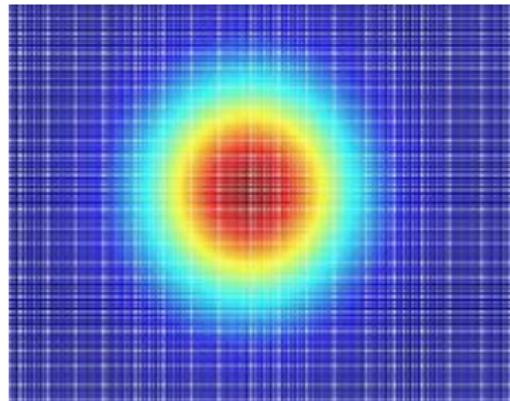
$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2} (x - \boldsymbol{\mu})^T \Sigma^{-1} (x - \boldsymbol{\mu})\right\},$$

where  $\boldsymbol{\mu}$  is the **mean vector** of  $X$ ,  $\Sigma$  is the **covariance matrix** of  $X$ .

# 1. Introduction

Independent case:

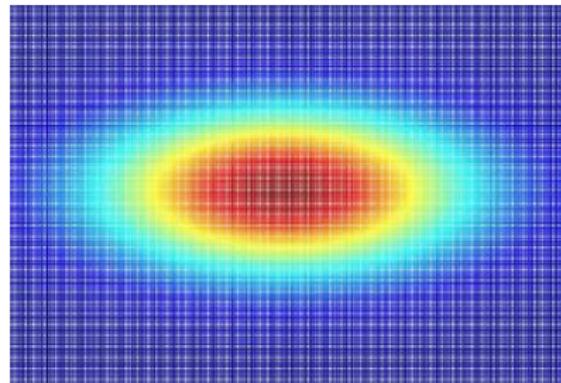
$$X = (X_1, \dots, X_n)^T, X_i \stackrel{iid}{\sim} N(0, \sigma^2).$$



# 1. Introduction

Heterogeneous independent case:

$$X = (X_1, \dots, X_n)^T, X_i \stackrel{iid}{\sim} N(0, \sigma_i^2).$$



# 1. Introduction

In general,  $\Sigma$  is a **positive definite matrix**. According to eigen decomposition:

$$\Sigma = \Gamma^T \Lambda \Gamma,$$

where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\Gamma^T \Gamma = I_n$ . Define the **square root of  $\Sigma$**  as:

$$\Sigma^{\frac{1}{2}} = \Gamma^T \Lambda^{\frac{1}{2}} \Gamma.$$

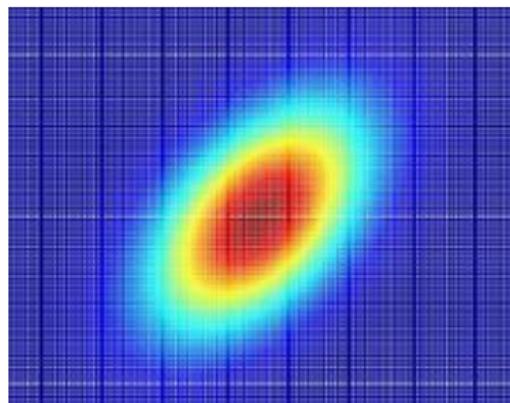
$\Sigma^{\frac{1}{2}}$  is also **positive definite**.

# 1. Introduction

Let  $\mathbf{Z} \sim N(0, I_n)$ , define:

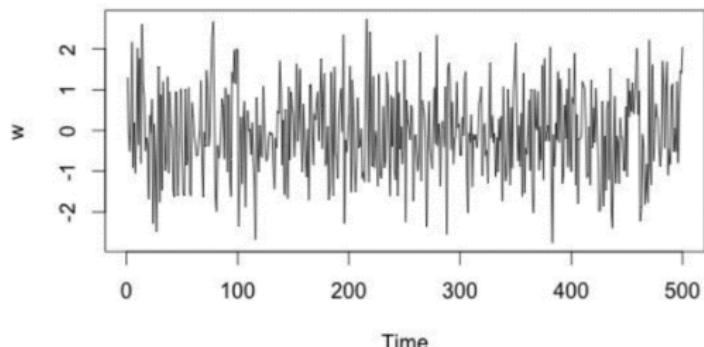
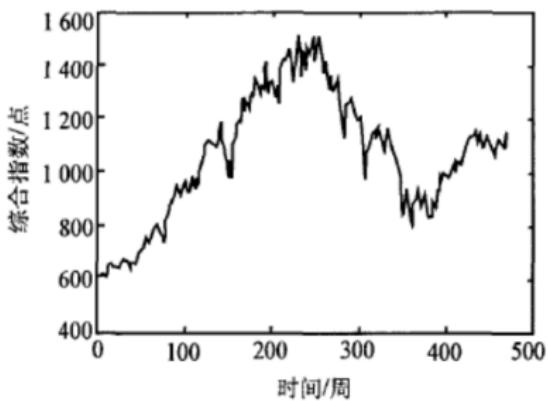
$$\mathbf{X} = \Sigma^{\frac{1}{2}} \mathbf{Z} + \boldsymbol{\mu},$$

then  $E(\mathbf{X}) = \boldsymbol{\mu}$ ,  $Cov(\mathbf{X}) = \Sigma$ .



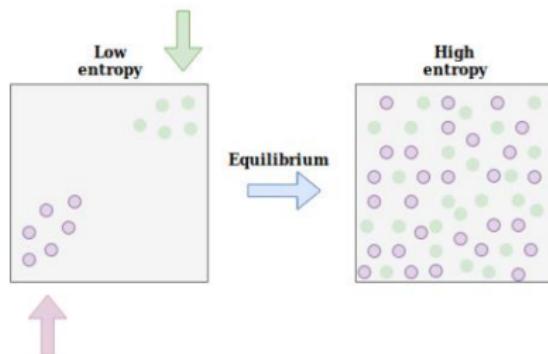
# 1. Introduction

Example:



## 2. Why Gaussian Assumption?

**Shannon Entropy:** Measure the level of chaos in a system.



## 2. Why Gaussian Assumption?

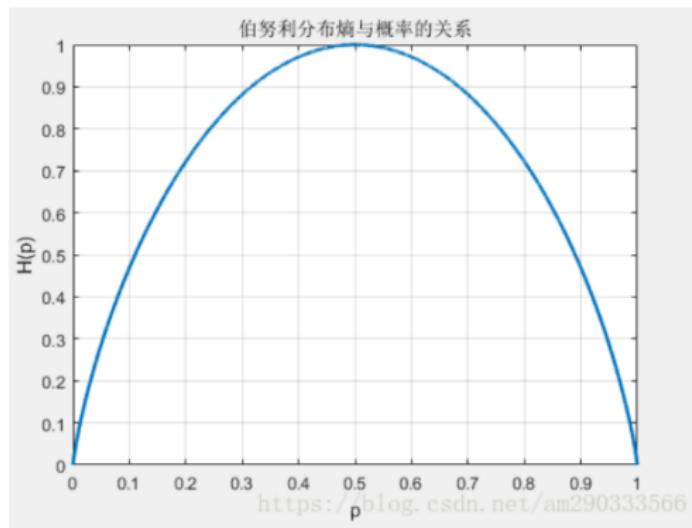
Given a random element  $X$  with pdf:  $p(x)$ , the entropy of  $X$  is defined as:

$$H(p) = - \int p(x) \log p(x) dx,$$

which is quantified the **randomness** of  $X$ .

## 2. Why Gaussian Assumption?

Example:



## 2. Why Gaussian Assumption?

### Theorem 4.1.2

Let  $q(x)$  be the pdf of random vector  $X$ , which has

$$E(X) = \mu, Cov(X) = \Sigma,$$

then  $H(q) \leq H(p)$ , where  $p \sim N(\mu, \Sigma)$ ,

Proof:

$$\begin{aligned} 0 &\leq KL(q || p) = \int_{-\infty}^{+\infty} q(x) \log \frac{q(x)}{p(x)} dx \\ &= -h(q) - \int q(x) \log p(x) dx \\ &= -h(q) - \int q(x) \left( \log \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} - \frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu) \right) dx \\ &= -h(q) + \frac{D}{2} \log 2\pi + \frac{1}{2} \log |\Sigma| + \frac{1}{2} \int q(x) (x-\mu)^T \Sigma^{-1} (x-\mu) dx \\ &= -h(q) + h(p) \\ \therefore h(p) &> h(q) \end{aligned}$$

## 2. Why Gaussian Assumption?

There are mainly two reason for the choice of Gaussian noise:

- **Maximize the randomness** of noise if the second moment of noise exists.
- Lindeberg-Feller Central limit theorem.

### 3. MLE for Gaussian

#### Theorem 4.2

$X_i \sim N(\mu, \Sigma)$   $i = 1, \dots, n$  then the MLE for  $\mu, \Sigma$  is given by

$$\hat{\mu} = \frac{1}{n} \sum X_i, \quad \hat{\Sigma} = \frac{1}{n} \sum (X_i - \bar{X})(X_i - \bar{X})^T$$

Proof : The Likelihood function:

$$x^T A x = \text{tr}(x^T A x) = \text{tr}(x x^T A) = \text{tr}(A x x^T)$$

$$\begin{aligned} L(\mu, \Sigma) &= \prod_{i=1}^n \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1} (x_i - \mu)) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)^T \Sigma^{-1} (x_i - \mu)\right) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left[-\frac{1}{2} \sum_{i=1}^n \text{tr}((x_i - \mu)^T \Sigma^{-1} (x_i - \mu))\right] \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left(-\frac{1}{2} \sum_{i=1}^n \text{tr}(\Sigma^{-1} (x_i - \mu)(x_i - \mu)^T)\right) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left(\text{tr}\left(-\frac{1}{2} \Sigma^{-1} \sum_{i=1}^n (x_i - \bar{x} + \bar{x} - \mu)(x_i - \bar{x} + \bar{x} - \mu)^T\right)\right) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left(\text{tr}\left(-\frac{1}{2} \Sigma^{-1} (A + n(\bar{x} - \mu)(\bar{x} - \mu)^T)\right)\right) \\ &= \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{n}{2}}} \exp\left[-\text{tr}\left(\frac{1}{2} \Sigma^{-1} A\right) - \frac{n}{2} (\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)^T\right] \\ &\Rightarrow \hat{\mu} = \bar{x} \end{aligned}$$

### 3. MLE for Gaussian

Lemma: Let  $B$  be  $m \times m$  positive definite matrix, then  
 $\text{tr}(B) - \ln|B| \geq m$   
the equality holds iff  $B = I_m$

证明

因为  $B > 0$ , 所以  $B$  的全部特征值  $\lambda_1, \dots, \lambda_m > 0$ , 且

$$|B| = \lambda_1 \cdots \lambda_m$$

利用不等式  $\ln(1+x) \leq x$  (当  $x+1>0$ ), 可得

$$\begin{aligned} \ln|B| &= \sum_{i=1}^m \ln \lambda_i = \sum_{i=1}^m \ln(1 + \lambda_i - 1) \\ &\leq \sum_{i=1}^m (\lambda_i - 1) = \text{tr}(B) - m \end{aligned}$$

所以

$$\text{tr}B - \ln|B| \geq m$$

当且仅当  $\lambda_i$  均等于 1 时等号成立, 即  $B = I_p$

### 3. MLE for Gaussian

$$\begin{aligned}\ln L(\bar{x}, \Sigma) &= -\frac{nm}{2} \ln(2\pi) - \frac{n}{2} \ln|\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} A) \\&= -\frac{nm}{2} \ln(2\pi) - \frac{n}{2} [\ln|\Sigma| + \text{tr}(\Sigma^{-1} \frac{A}{n})] \\&= C - \frac{n}{2} [\text{tr}(\Sigma^{-1} \frac{A}{n}) - \ln|\Sigma^{-1} \frac{A}{n}| + \ln|\frac{A}{n}|] \\&= C - \frac{n}{2} [\text{tr}(\Sigma^{-\frac{1}{2}} \frac{A}{n} \Sigma^{-\frac{1}{2}}) - \ln|\Sigma^{-\frac{1}{2}} \frac{A}{n} \Sigma^{-\frac{1}{2}}| + \ln|\frac{A}{n}|] \\&\Rightarrow \hat{\Sigma} = \frac{A}{n} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top\end{aligned}$$

## Appendix: Information theory

KL divergence:

$$KL(q||p) = \int p(x) \log \frac{p(x)}{q(x)} dx.$$

Conditional Entropy:

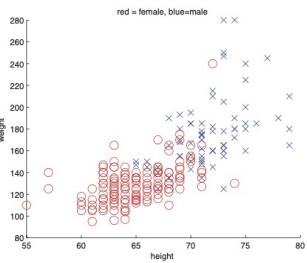
$$H(Y|X) = - \iint p(x, y) \log p(y|x) dx dy.$$

Joint Entropy:

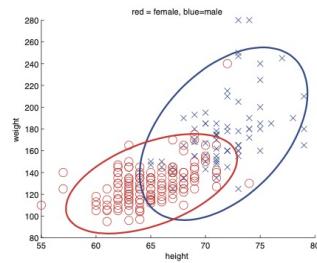
$$H(X, Y) = - \iint p(x, y) \log p(x, y) dx dy.$$

## Part II: Gaussian discriminant analysis

Classification problem: Given  $X$ , we want to know  $Y=0$  or  $Y=1$ .



(a)



(b)

discriminative learning      v.s.      generative Learning

Find  $P(Y|X)$

$$f: X \rightarrow \{0, 1\}$$

Find joint  $P(X, Y)$ ,  $P(X=0|Y)$ ,  $P(X=1|Y)$   
priors  $P(Y)$   $\Rightarrow P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$

Remark:  $\frac{P(Y=0|X)}{P(Y=1|X)} = \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}$  Bayes classifier.

• Linear discriminant analysis

Assumption:  $X|Y=0 \sim N(\mu_0, \Sigma)$  same covariance matrix

$X|Y=1 \sim N(\mu_1, \Sigma) \quad Y \sim Ber(p)$

$$P(Y=c|X) \propto \pi_c \exp[\mu_c^\top \Sigma^{-1} X - \frac{1}{2} X^\top \Sigma^{-1} X - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c]$$

$$\propto \exp[\mu_c^\top \Sigma^{-1} X - \frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c + \log \pi_c] \exp[-\frac{1}{2} X^\top \Sigma^{-1} X]$$

$$\text{Let } \beta_c = \Sigma^{-1} \mu_c \quad r_c = -\frac{1}{2} \mu_c^\top \Sigma^{-1} \mu_c + \log \pi_c$$

$$P(Y=c|X) = \frac{\exp(\beta_c^\top X + r_c)}{\exp(\beta_0^\top X + r_0) + \exp(\beta_1^\top X + r_1)} = \text{sigm}((\beta_c - \beta_0)^\top X + (r_c - r_0))$$

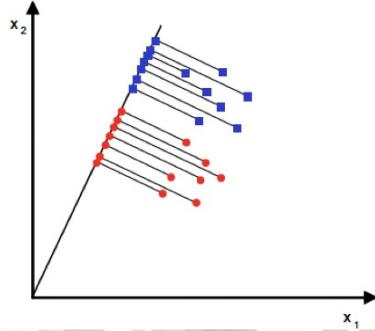
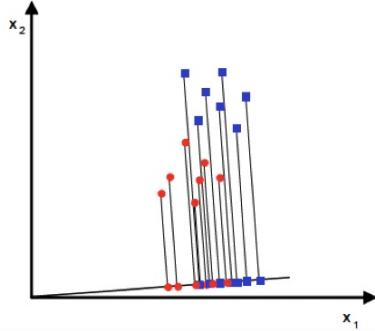
$$\text{Let } w = \beta_1 - \beta_0, \quad x_0 = \frac{1}{2}(\mu_1 + \mu_0) - (\mu_1 - \mu_0) \frac{\log(\pi_1/\pi_0)}{(\mu_1 - \mu_0)^\top \Sigma^{-1} (\mu_1 - \mu_0)}$$

$$P(Y=c|X) = \text{sigm}(w^\top (X - x_0)), \text{ if } w^\top (X - x_0) \Rightarrow Y=1$$

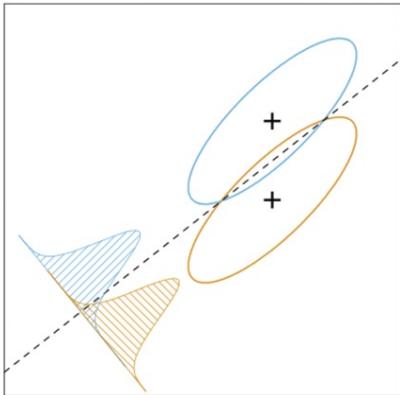
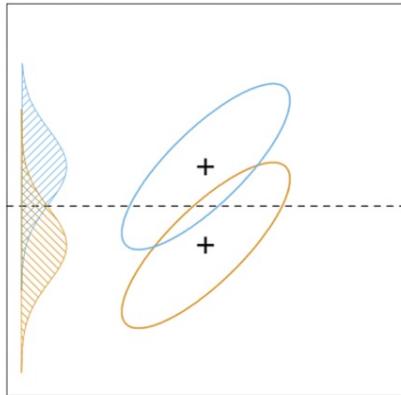
$$\text{Remark: } w^\top (X - x_0) > 0 \iff \log \frac{P(X|Y=1)}{P(X|Y=0)} > 0 \iff \frac{P(X|Y=1)}{P(X|Y=0)} > 1$$

# Projection and LDA

- separable



- inseparable



LDA is equal to:

$$\max_{\alpha} \frac{(E(\alpha^T x_1) - E(\alpha^T x_2))^2}{\text{var}(\alpha^T x_1 - \alpha^T x_2)} \Rightarrow \alpha \propto \Sigma^{-1} (\mu_1 - \mu_2)$$

## • Quadratic discriminant analysis

Assumption:  $X|Y=0 \sim N(\mu_0, \Sigma_0)$   $\Sigma_0 \neq \Sigma_1$

$X|Y=1 \sim N(\mu_1, \Sigma_1)$

$$P(Y=c|X) = \frac{\pi_c / 2\pi |\Sigma_c|^{-\frac{1}{2}} \exp[-\frac{1}{2}(X-\mu_c)^T \Sigma_c^{-1} (X-\mu_c)]}{\sum_{c=0}^1 \pi_c / 2\pi |\Sigma_c|^{-\frac{1}{2}} \exp[-\frac{1}{2}(X-\mu_c)^T \Sigma_c^{-1} (X-\mu_c)]}$$

$$\log P(Y=c|X) \propto -\frac{1}{2} \log |\Sigma_c| - \frac{1}{2} (X-\mu_c)^T \Sigma_c^{-1} (X-\mu_c) + \log \pi_c$$

$$\log \frac{P(Y=1|X)}{P(Y=0|X)} = \underbrace{\left\{ \frac{1}{2} \log \frac{|\Sigma_1|}{|\Sigma_0|} + \mu_1^T \Sigma_1^{-1} \mu_1 - \mu_0^T \Sigma_0^{-1} \mu_0 \right\}}_{\beta_0} + \underbrace{\left\{ \Sigma_1^{-1} \mu_1 - \Sigma_0^{-1} \mu_0 \right\}^T X - \frac{1}{2} X^T (\Sigma_1^{-1} - \Sigma_0^{-1}) X}_{\beta_1} + \underbrace{\beta_2}_0$$

$$\text{If } \beta_0 + \beta_1^T X + X^T \beta_2 X > 0 \implies Y=1$$

Remark: In LDA, discriminant function is linear function.

In QDA, discriminant function is quadratic function.

## • How to estimate parameters in discriminant analysis?

Solution: use MLE

$$\begin{cases} \hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} x_i \\ \hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T \end{cases}$$

In high dimensions, the MLE may badly overfit. ( $N_c < D$ )

Solution: ① Assume  $\Sigma_c$  is diagonal  $\implies$  naive Bayes

② Assume  $\Sigma_c = \Sigma$ .

③ Projection.

④ MAP estimation.

## • Diagonal LDA

In RDA, let  $\lambda = \frac{b}{N} \Rightarrow$  diagonal LDA

$$\log P(Y=c | X) = -\sum_{j=1}^c \frac{(x_{ij} - \mu_{cj})^2}{2\hat{\sigma}_j^2} + \log \pi_c$$

$$\mu_{cj} = \bar{x}_{cj}, \hat{\sigma}_j^2 = \frac{\sum_{i=1}^N (x_{ij} - \bar{x}_{cj})^2}{N-c}$$

**Remark:** diagonal LDA is a special case of naive Bayes.

Naive Bayes classifier

Assumption: Give a class  $Y=j$ , features  $X_k$  are independent

$$\begin{aligned} \log \frac{P(Y=j | X)}{P(Y=j | X)} &= \log \frac{\pi_j P(X | Y=j)}{\pi_j P(X | Y=j)} \\ &= \log \frac{\pi_j \prod_{k=1}^p P(X_k | Y=j)}{\pi_j \prod_{k=1}^p P(X_k | Y=j)} \\ &= \log \frac{\pi_j}{\pi_j} + \sum_{k=1}^p \log \left( \frac{P(X_k | Y=j)}{P(X_k | Y=j)} \right) \\ &= \alpha_j + \sum_{k=1}^p \alpha_{jk}(X_k) \quad \text{generalized additive model} \end{aligned}$$

## • Nearest shrunken centroids classifier (NSCC)

Motivation: diagonal LDA depends on all of the features.

In high-dimension problem, we only want to a subset.

We let some  $\underline{x}_{cj} = \bar{x}_j$ , then drop gene  $j$  in classification.

Method:

① Shrinking the classwise mean toward the overall mean.

$$d_{cj} = \frac{\underline{x}_{cj} - \bar{x}_j}{m_c(s_j + s_0)}$$

where  $m_c^2 = \frac{1}{N_c} - \frac{1}{N}$ ,  $s_0 > 0$  is a small constant median( $s_j$ )

② Use soft thresholding

$$d_{cj} = \text{sign}(d_{cj})(|d_{cj}| - \Delta) +$$

$\Delta$  is a parameter to be determined (may use CV)

Remark: hard thresholding:  $d'_{cj} = d_{cj} I(|d_{cj}| > \Delta)$

③ obtain the shrunken version of  $\bar{x}_{c_j}'$

$$\bar{x}_{c_j}' = \bar{x}_j + m_k(s_j + s_0)d_{c_j}'$$

Remark:  $\bar{x}_{c_j}'$  is also a Lasso-style estimator.

④ use  $\bar{x}_{c_j}'$  in discriminant score

$$\log P(Y=c|X) = -\sum_{j=1}^D \frac{(x_j - \mu_{c_j})^2}{2\hat{\sigma}_j^2} + \log \pi_c$$

$$\hat{\mu}_{c_j} = \bar{x}_{c_j}', \quad \hat{\sigma}_j^2 = \frac{\sum_{i:y_i=c}(x_{ij} - \bar{x}_{c_j}')^2}{N-c}$$

Remark: If  $d_{c_j}' = 0$ , then  $\bar{x}_{c_j}' = \bar{x}_j$ , gene j don't play a role in classification.

- LDA vs logistic regression

$$\text{LDA: } \log \frac{P(Y=1|X)}{P(Y=0|X)} = \log \frac{\pi_1}{\pi_0} - \frac{1}{2}(\mu_1 + \mu_0)^T \Sigma^{-1}(\mu_1 - \mu_0) + x^T \Sigma^{-1}(\mu_1 - \mu_0)$$

$$= \alpha_0 + \alpha_1^T X$$

$$\text{logistic: } \log \frac{P(Y=1|X)}{P(Y=0|X)} = \beta_0 + \beta_1^T X$$

① LDA need specify joint distribution, but logistic not.

②  $P(X, Y=k) = P(X) P(Y=k|X)$

Both LDA and logistic have the logit-linear form.

③ logistic treats  $P(X)$  as an arbitrary density function.

$$\text{LDA } P(X, Y=k) = N(X|\mu_k, \Sigma) \pi_k$$

$$P(X) = \sum_{k=0}^1 N(X|\mu_k, \Sigma) \pi_k$$

④ Logistic fits  $P(Y|X)$  by maximizing conditional likelihood.

LDA fits  $P(Y|X)$  by MLE.

⑤ logistic: outlier are down-weighted.

LDA: no robust.

### Part III: Evaluation in jointly Gaussian distributions

Given a joint distribution  $P(X_1, X_2)$ , we consider marginals  $P(X_1)$  and conditionals  $P(X_1 | X_2)$ .

Thm 4.3.1 Suppose  $x = (X_1, X_2)$  is jointly Gaussian with parameters  
 $\mu = (\mu_1 \ \mu_2)$   $\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$ .  $\Lambda = \Sigma^{-1} = \begin{pmatrix} \Lambda_{11} & \Lambda_{12} \\ \Lambda_{21} & \Lambda_{22} \end{pmatrix}$

Then,

$$P(X_1) = N(X_1 | \mu_1, \Sigma_{11})$$

$$P(X_2) = N(X_2 | \mu_2, \Sigma_{22})$$

$$P(X_1 | X_2) = N(X_1 | \mu_{1|2}, \Sigma_{1|2})$$

$$\text{where } \mu_{1|2} = \Sigma_{12}(\Lambda_{11}\mu_1 - \Lambda_{12}(X_2 - \mu_2))$$

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

Proof : ① linear transform

$$\mathbf{Z} = \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} = \begin{pmatrix} I_p - \Sigma_{12}\Sigma_{22}^{-1} \\ 0 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \begin{pmatrix} X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 \\ X_2 \end{pmatrix} = BX$$

$Z$  is MVN.

$$E(Z) = BE(X) = \begin{pmatrix} \mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2 \\ \mu_2 \end{pmatrix}$$

$$V(Z) = B V(X) B^T = \begin{pmatrix} \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} & 0 \\ 0 & \Sigma_{22} \end{pmatrix}$$

$$\text{② } Z \sim N(\mu_1 - \Sigma_{12}\Sigma_{22}^{-1}\mu_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}) N(\mu_2, \Sigma_{22})$$

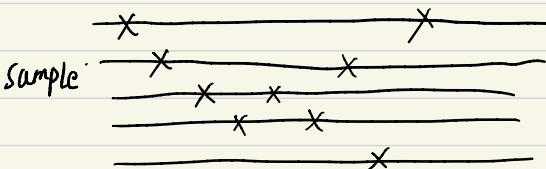
$$\left\{ \begin{array}{l} Z_1 = X_1 - \Sigma_{12}\Sigma_{22}^{-1}X_2 \\ Z_2 = X_2 \end{array} \right.$$

$$X_1 | X_2 \sim N(\mu_1 - \Sigma_{12}\Sigma_{22}^{-1}(X_2 - \mu_2), \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

## Some example

- Data imputation

Aim: Predict the missing entries  
feature



① Compute  $P(X_{hi} | X_{vi}, \theta)$ ,  $hi$ : hidden,  $vi$ : visible.

② Compute the marginal:  $P(X_{hi} | X_{vi}, \theta)$

③  $\hat{X}_{ij} = E[X_{ij} | X_{vi}, \theta]$

## Part IV : Linear Gaussian systems

$X \in \mathbb{R}^{D_X}$  hidden variable ,  $Y \in \mathbb{R}^{D_Y}$  noise observation of  $X$

Assumption:  $P(X) = N(X | \mu_X, \Sigma_X)$

$$P(Y|X) = N(Y|AX+b, \Sigma_Y)$$

Aim: infer  $X$  from  $Y$ .

$$\text{Thm 4.4.1 } P(X|Y) = N(X | \mu_{X|Y}, \Sigma_{X|Y}) , \Sigma_{X|Y}^{-1} = \Sigma_X^{-1} + A^T \Sigma_Y^{-1} A$$

$$\mu_{X|Y} = \Sigma_{X|Y} [A^T \Sigma_Y^{-1} (Y - b) + \Sigma_X^{-1} \mu_X]$$

$$P(Y) = N(Y | A\mu_X + b, \Sigma_Y + A\Sigma_X A^T)$$

Proof: ① compute  $P(X|Y) = P(X)P(Y|X)$

$$\log P(Y|X) \propto -\frac{1}{2}(X - \mu_X)^T \Sigma_X^{-1} (X - \mu_X) - \frac{1}{2}(Y - AX - b)^T \Sigma_Y^{-1} (Y - AX - b) \Rightarrow \text{gaussian}$$

Ignore linear and constant terms:

$$\begin{aligned} Q &= -\frac{1}{2} X^T \Sigma_X^{-1} X - \frac{1}{2} Y^T \Sigma_Y^{-1} Y - \frac{1}{2} (AX)^T \Sigma_Y^{-1} (AX) + Y^T \Sigma_Y^{-1} AX \\ &= -\frac{1}{2} (X)^T \begin{pmatrix} \Sigma_X^{-1} + A^T \Sigma_Y^{-1} A & -A^T \Sigma_Y^{-1} \\ -\Sigma_Y^{-1} A & \Sigma_Y^{-1} \end{pmatrix} (X) \\ &= -\frac{1}{2} (X)^T \Sigma^{-1} (X) \end{aligned}$$

② use Thm 4.3.1 , let  $\Sigma^{-1} = \begin{pmatrix} \lambda_{XX} & \lambda_{XY} \\ \lambda_{YX} & \lambda_{YY} \end{pmatrix}$

$$P(X|Y) = N(\mu_{X|Y}, \Sigma_{X|Y})$$

$$\Sigma_{X|Y} = \lambda_{XX}^{-1} = (\Sigma_X^{-1} + A^T \Sigma_Y^{-1} A)^{-1}$$

$$\mu_{X|Y} = \Sigma_{X|Y} L \Sigma_X^{-1} \mu_X + A^T \Sigma_Y^{-1} (Y - b))$$

Example:  $\mu \rightarrow X \rightarrow Y$

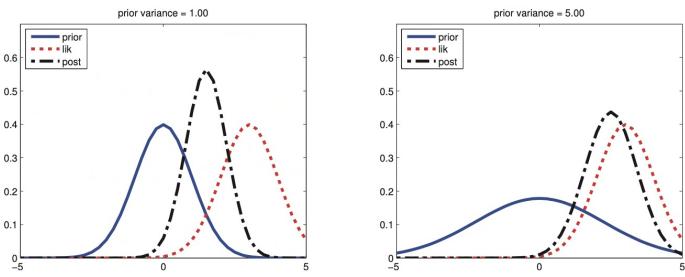
- Inferring an unknown scalar from noisy measurements.

Assumption:  $P(Y_i|X) = N(Y_i | X, \lambda_Y^{-1}) , \lambda_Y = \frac{1}{\delta^2}$

$$P(X) = N(X | \mu_0, \lambda_0^{-1})$$

$$P(X|Y) = N(X | \mu_N, \lambda_N^{-1})$$

$$\text{where } \lambda_N = \lambda_0 + N\lambda_Y \quad \mu_N = \frac{N\lambda_Y \bar{Y} + \lambda_0 \mu_0}{\lambda_N} = \frac{N\lambda_Y}{N\lambda_Y + \lambda_0} \bar{Y} + \frac{\lambda_0}{N\lambda_Y + \lambda_0} \mu_0 \quad (\text{convex combination})$$



**Figure 4.12** Inference about  $x$  given a noisy observation  $y = 3$ . (a) Strong prior  $\mathcal{N}(0, 1)$ . The posterior mean is "shrunk" towards the prior mean, which is 0. (a) Weak prior  $\mathcal{N}(0, 5)$ . The posterior mean is similar to the MLE. Figure generated by gaussInferParamsMeanid.

$$P(x|D, \sigma^2) = \mathcal{N}(x | \mu_N, T_N^{-2})$$

$$\begin{aligned} T_N^{-2} &= \frac{1}{\lambda_N} = \frac{1}{\lambda_0 + N\lambda_y} = \frac{1}{\frac{N\lambda_y}{N\lambda_0 + \lambda_y} + \frac{\lambda_0}{\lambda_0 + \lambda_y}} = \frac{\sigma^2 T_0^{-2}}{N T_0^{-2} + \sigma^2}, \text{ where } T_0^{-2} = \frac{1}{\lambda_0} \\ \mu_N &= \frac{1}{N T_0^{-2} + \sigma^2} \mu_0 + \frac{N \lambda_y}{N T_0^{-2} + \sigma^2} y \end{aligned}$$

Special case  $N=1$ : (define  $\Sigma_y = \sigma^2$ ,  $\Sigma_0 = T_0^{-2}$ ,  $\Sigma_1 = T_1^{-2}$ )

$$P(x|y) = \mathcal{N}(x | \mu_1, \Sigma_1)$$

$$\Sigma_1 = \frac{\Sigma_y \Sigma_0}{\Sigma_0 + \Sigma_y} = \left( \frac{1}{\Sigma_0} + \frac{1}{\Sigma_y} \right)^{-1}$$

$$\begin{aligned} \mu_1 &= \Sigma_1 \left( \frac{\mu_0}{\Sigma_0} + \frac{y}{\Sigma_y} \right) \\ &= \frac{\Sigma_y}{\Sigma_y + \Sigma_0} \mu_0 + \frac{\Sigma_0}{\Sigma_y + \Sigma_0} y \\ &= \mu_0 + (y - \mu_0) \frac{\Sigma_0}{\Sigma_y + \Sigma_0} \\ &= y - (y - \mu_0) \frac{\Sigma_y}{\Sigma_y + \Sigma_0} \end{aligned}$$

Convex combination

Prior mean adjusted toward data

Shrinkage

Signal-to-noise ratio  $\frac{\mathbb{E} x^2}{\mathbb{E} \Sigma^2} = \frac{\Sigma_0 + \mu_0^2}{\Sigma_y}$

- Inferring an unknown vector from noisy measurements

observation:  $y_i \sim \mathcal{N}(x, \Sigma_y)$ , Gaussian prior:  $x \sim \mathcal{N}(\mu_0, \Sigma_0)$

Aim: Find  $x$ . (Y may has multiple observations with different  $\Sigma$ )

$$P(x|y_1, \dots, y_N) = \mathcal{N}(x | \mu_N, \Sigma_N)$$

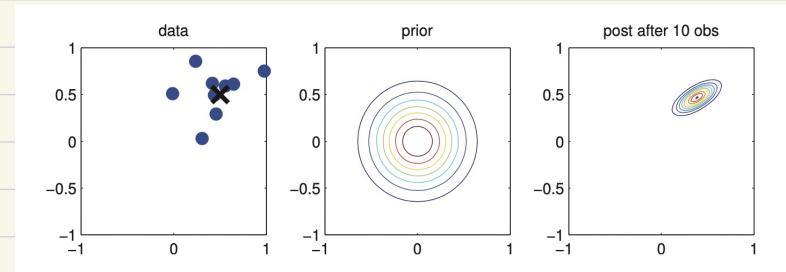
$$\mu_N = \Sigma_N (\Sigma_y^{-1} (N \bar{y}) + \Sigma_0^{-1} \mu_0)$$

$$\Sigma_N^{-1} = \Sigma_0^{-1} + N \Sigma_y^{-1}$$

Simple case:

$$Y_i \sim N(X, \Sigma_Y)$$

$$\text{Prior: } N(X(0, 0.1I))$$



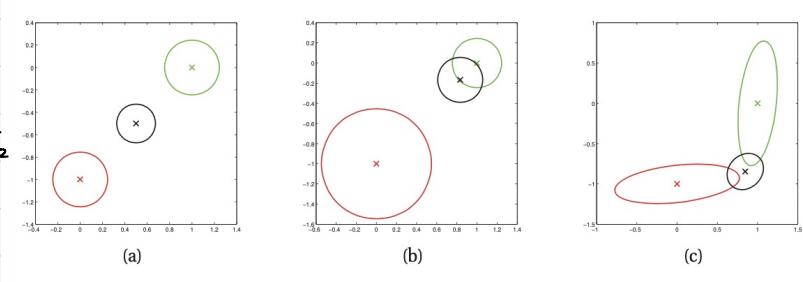
$$Y_1 = (0, -1), Y_2 = (1, 0)$$

$$(a) \Sigma_{Y_1} = \Sigma_{Y_2} = 0.01 I_2$$

$$(b) \Sigma_{Y_1} = 0.05 I_2 \quad \Sigma_{Y_2} = 0.01 I_2$$

$$(c) \Sigma_{Y_1} = 0.01 \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$\Sigma_{Y_2} = 0.01 \begin{pmatrix} 1 & 1 \\ 1 & 10 \end{pmatrix}$$

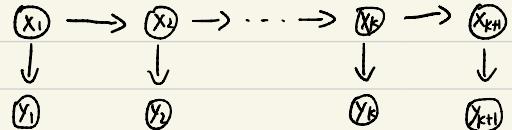


• An example : Kalman Filter

① Model:

$$X_k = F X_{k-1} + B U_{k-1} + W_{k-1}$$

$$Y_k = H X_k + V_k$$



where  $F$  is the state transition matrix,  
 $U_k$  is control vector,

$$W_{k-1} \sim N(0, Q) \text{ noise}$$

$H$  is measurement matrix.

$$V_k \sim N(0, R) \text{ noise}$$

We need to estimate  $X_k$  given  $X_0, Y_1, \dots, Y_k$

## ② Algorithm

$$E(x_n | y_1, \dots, y_n) \xrightarrow{\text{prediction}} E(x_{n+1} | y_1, \dots, y_n) \xrightarrow{\text{correction}} E(x_{n+1} | y_1, \dots, y_{n+1})$$

Prediction:

- (a) state estimate  $\hat{x}_k^- = F \hat{x}_{k-1}^+ + B u_{k-1}$
- (b) error covariance  $P_k^- = F P_{k-1}^+ F^T + Q$

Update:

- (a) measurement residual:  $\tilde{y}_k = y_k - H \hat{x}_k^-$
- (b) Kalman gain:  $K_k = P_k^- H^T (R + H P_k^- H^T)^{-1}$
- (c) state estimate:  $\hat{x}_k^+ = \hat{x}_k^- + K_k \tilde{y}_k$
- (d) error covariance:  $P_k^+ = (I - K_k H) P_k^-$   
(-: predicted (prior) estimate; +: updated (posterior) estimate)

## ③ Application

A robot.

$$X_k = [P, V] \quad P: \text{position}, V: \text{velocity}$$

$$F = \begin{bmatrix} 1 & \Delta t \\ 0 & 1 \end{bmatrix} \quad P_k = P_{k-1} + \Delta t V_{k-1}$$

$$V_k = V_{k-1}$$

$U_k$ : acceleration

$$B = \begin{bmatrix} \frac{\Delta t^2}{2} \\ \Delta t \end{bmatrix} \quad P_k = P_{k-1} + \Delta t V_{k-1} + \frac{1}{2} a \Delta t^2$$

$$V_k = V_{k-1} + a \Delta t$$

$$P_k = F_k P_{k-1} F_k^T + Q_k$$