# CS5014 Machine Learning
## Lecture 13 Unsupervised Learning

Lei Fang

School of Computer Science, University of St Andrews

Spring 2021

# Some responses: geometry of ridge regression

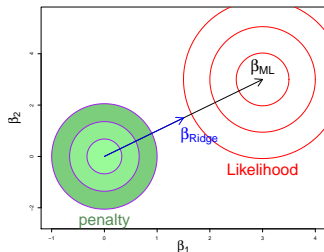$$\beta_{\text{ridge}} \equiv \underset{\beta}{\text{argmin}} \underbrace{||\boldsymbol{y} - \boldsymbol{X}\beta||_2^2}_{L(\theta)} + \lambda||\beta||_2^2$$

$$\beta_{\text{ridge}} = (\boldsymbol{X}^T\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}\boldsymbol{y}$$



- assume feature vectors of $\boldsymbol{X}$ are norm 1 and orthogonal i.e. orthogonormal: $\boldsymbol{X}^T\boldsymbol{X} = \boldsymbol{I}$

$$\beta_{\text{ridge}} = \frac{1}{1+\lambda}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}\boldsymbol{y} = \frac{1}{1+\lambda}\beta_{ML}$$
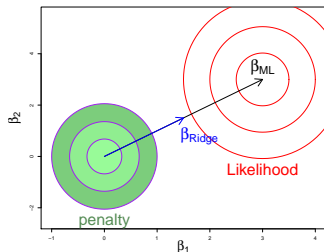
- why the likelihood (red) contours are circular for this case ?

# Some responses: geometry of ridge regression

$$\beta_{\text{ridge}} \equiv \underset{\beta}{\text{argmin}} \underbrace{||\boldsymbol{y} - \boldsymbol{X}\beta||_2^2}_{L(\boldsymbol{\theta})} + \lambda||\beta||_2^2$$

$$\beta_{\text{ridge}} = (\boldsymbol{X}^T\boldsymbol{X} + \lambda\boldsymbol{I})^{-1}\boldsymbol{X}\boldsymbol{y}$$



- assume feature vectors of $\boldsymbol{X}$ are norm 1 and orthogonal i.e. orthogonormal: $\boldsymbol{X}^T\boldsymbol{X} = \boldsymbol{I}$

$$\beta_{\text{ridge}} = \frac{1}{1+\lambda}(\boldsymbol{X}^T\boldsymbol{X})^{-1}\boldsymbol{X}\boldsymbol{y} = \frac{1}{1+\lambda}\beta_{ML}$$

- why the likelihood (red) contours are circular for this case ? remember the Hessian of $L(\boldsymbol{\theta})$?
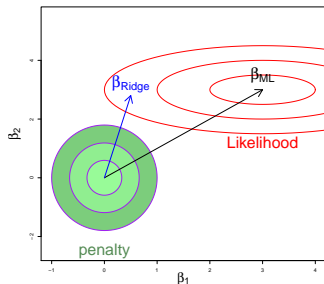
$$H_L = 2\boldsymbol{X}^T\boldsymbol{X} \propto \boldsymbol{I}$$

- uniform discount: $\beta_k$ receives the same discount: $\frac{1}{1+\lambda}$

# Some responses: some other cases

When $H_L \neq I$, the shrinkage is **NOT** uniform;

$$\beta_k^{\text{ridge}} = \frac{\sigma_k}{\sigma_k + \lambda} \beta_k^{ML}$$



- $\sigma_k$ is the directional curvature of $X^T X$ (also the eigen value)
- flat curve ($\beta_1$ direction) $\Rightarrow \sigma_k$ smaller $\Rightarrow$ more discount
- curvy curve ($\beta_2$ direction) $\Rightarrow \sigma_k$ larger $\Rightarrow$ less discount
- makes perfect sense!
  - flat means less confident (or large variance): shrink more
  - peak means confident estimate (small variance): shrink less

($*$) the equation is true when $X^T X$ is a diagonal matrix; up to basis translation for more general cases

# Today's topic

Unsupervised learning
- clustering
- k-means

Revisit multivariate Gaussian

Revisit k-means
- mixture of Gaussians
- EM algorithm for mixture model
  - K-means is just a specific case
- other kinds of mixture models

# Unsupervised learning

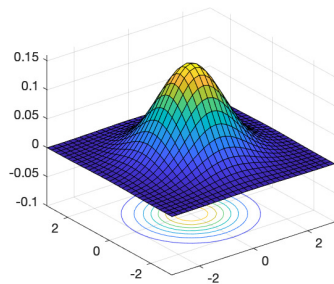# Clustering

# K-means

# K-means

# Demonstration

# Limitations of K-means

# Dissect multivariate Gaussians

$$p(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}\underbrace{(\boldsymbol{x}-\boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}_{d_{\boldsymbol{\Sigma}}}\right]$$

- a distance measure : (aka mahalanobis distance)

  $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu})$ : between $\boldsymbol{x}$ and $\boldsymbol{\mu}$
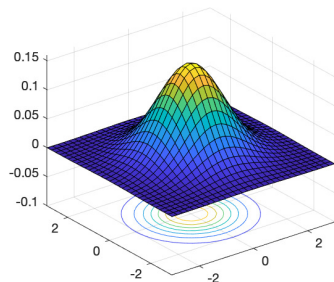
# Dissect multivariate Gaussians

$$p(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[ -\frac{1}{2} \underbrace{(\boldsymbol{x} - \boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})}_{d_{\boldsymbol{\Sigma}}} \right]$$

- a distance measure : (aka mahalanobis distance)

  $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu})$ : between $\boldsymbol{x}$ and $\boldsymbol{\mu}$

- $-$ : $p$ is **negatively** related to the distance

  larger $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu}) \Rightarrow$ further away $\boldsymbol{x}$ from $\boldsymbol{\mu} \Rightarrow$ smaller $p$

# Dissect multivariate Gaussians

$$p(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[\; -\frac{1}{2} \underbrace{(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu})}_{d_{\boldsymbol{\Sigma}}} \right]$$
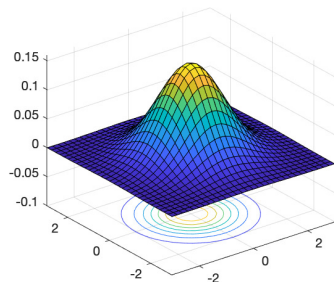
- a distance measure : (aka mahalanobis distance)

  $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu})$ : between $\boldsymbol{x}$ and $\boldsymbol{\mu}$

- $-$ : $p$ is **negatively** related to the distance

  larger $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu}) \Rightarrow$ further away $\boldsymbol{x}$ from $\boldsymbol{\mu} \Rightarrow$ smaller $p$

- exp: makes sure $p(\boldsymbol{x}) > 0$

# Dissect multivariate Gaussians

$$p(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \equiv \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left[ -\frac{1}{2} \underbrace{(\boldsymbol{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{x} - \boldsymbol{\mu})}_{d_{\boldsymbol{\Sigma}}} \right]$$

- a distance measure : (aka mahalanobis distance)

  $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu})$ : between $\boldsymbol{x}$ and $\boldsymbol{\mu}$
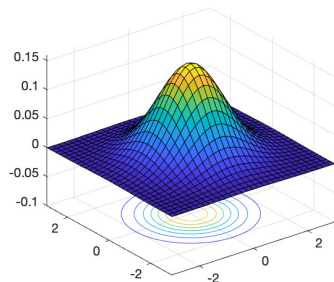
- $-$ : $p$ is **negatively** related to the distance

  larger $d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \boldsymbol{\mu}) \Rightarrow$ further away $\boldsymbol{x}$ from $\boldsymbol{\mu} \Rightarrow$ smaller $p$

- exp:  makes sure $p(\boldsymbol{x}) > 0$

- normalising constant: s.t. $\int N(\boldsymbol{x}; \cdot, \cdot) d\boldsymbol{x} = 1$;

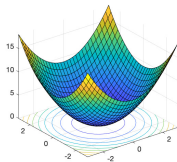  $|\boldsymbol{\Sigma}|$ : determinant; a volume measure-ish quantity
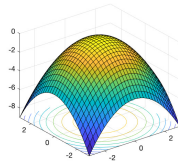
# Dissect multivariate Gaussians

**Key message**: $d_\Sigma$ (the distance) determines equal $p(\boldsymbol{x})$ levels

$$p(\boldsymbol{x}) \equiv \underbrace{\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}}}_{C:\text{normalising cst.}} \exp \Big[ \overbrace{- \frac{1}{2} \underbrace{(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}_{f_1 = d_\Sigma(\boldsymbol{x};\boldsymbol{\mu})}}^{f_2} \Big]$$

($f_3$ spans from $\exp$ to the end of the bracket.)



1. a distance measure:
   $f_1(\boldsymbol{x}) = d_\Sigma(\boldsymbol{x};\boldsymbol{\mu})$

2. negated distance:
   $f_2(\boldsymbol{x}) = -\frac{1}{2}f_1(\boldsymbol{x})$

3. exp. to make sure $p > 0$:
   $f_3(\boldsymbol{x}) = e^{f_2(\boldsymbol{x})}$

4. scaled to make sure $\int p(\boldsymbol{x})d\boldsymbol{x} = 1$:
   $p(\boldsymbol{x}) = C \cdot f_3(\boldsymbol{x})$

# Covariance matrix and distance

$$\mathbf{\Sigma} : \text{variance-covariance matrix}$$

- $d \times d$ symmetric matrix:

$$\mathbf{\Sigma} = \mathbf{\Sigma}^T$$

- positive definite (P.D.):

$$\mathbf{v}^T \mathbf{\Sigma} \mathbf{v} > 0, \quad \forall \mathbf{v} \in R^d$$

- why P.D. ? distance has to be positive ! (similar to univariate Gaussian: $(x - \mu)^2 \cdot \sigma^{-2} > 0$)

$$(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) > 0, \quad \text{where } \mathbf{v} = \mathbf{x} - \boldsymbol{\mu}$$

  - if $\mathbf{\Sigma}$ is P.D., then $\mathbf{\Sigma}^{-1}$ is also P.D.; so the above is a valid distance metric

Proof: Let $\mathbf{y} = \mathbf{\Sigma} \mathbf{v}$; then $\mathbf{y}^T \mathbf{\Sigma}^{-1} \mathbf{y} = \mathbf{v}^T \mathbf{\Sigma}^T \mathbf{\Sigma}^{-1} \mathbf{\Sigma} \mathbf{v} = \mathbf{v}^T \mathbf{\Sigma}^T \mathbf{v} = \mathbf{v}^T \mathbf{\Sigma} \mathbf{v} > 0$

# Diagonal $\mathbf{\Sigma}$: implies independence

If

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 & \ldots & 0 \\ 0 & \sigma_2^2 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \sigma_d^2 \end{bmatrix}; \quad \mathbf{\Sigma}^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 & \ldots & 0 \\ 0 & \frac{1}{\sigma_2^2} & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \frac{1}{\sigma_d^2} \end{bmatrix}$$

Then

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}|\mathbf{\Sigma}|^{1/2}} \exp\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T\mathbf{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\} = \frac{1}{(2\pi)^{d/2}(\prod_{i=1}^d \sigma_i^2)^{1/2}} \exp\{-\frac{1}{2}\sum_{i=1}^d (x_i-\mu_i)^2/\sigma_i^2\}$$

$$= \prod_{i=1}^d \underbrace{\frac{1}{(2\pi)^{1/2}\sigma_i} \exp\{-\frac{1}{2}(x_i-\mu_i)^2/\sigma_i^2\}}_{\text{unvariate Gaussian}} = \underbrace{\prod_{i=1}^d p(x_i)}_{\text{independence !}}$$
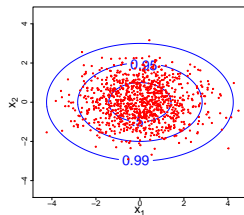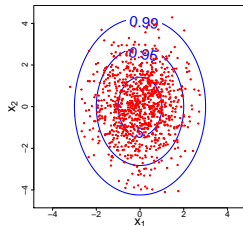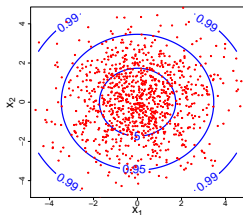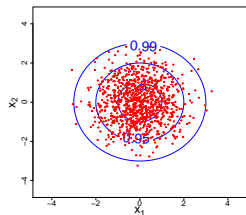
and each $p(x_i) = N(x_i; \mu_i, \sigma_i^2)$ is a univariate Gaussian

Remember independence ? it means knowing one does not inform the other: $p(x_i|\mathbf{x}_{/i}) = p(x_i)$

# Diagonal $\Sigma$: axis aligned ellipses

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} ; \Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix}$$ so $d_{\Sigma}(x; 0)$ are axis aligned ellipses
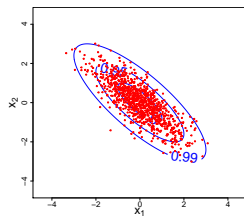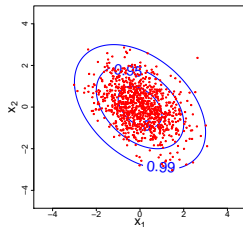


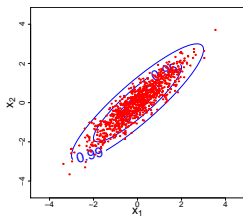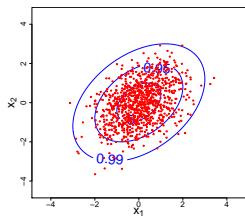$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \qquad \Sigma = \begin{bmatrix} 3 & 0 \\ 0 & 3 \end{bmatrix} \qquad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \qquad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$

University of St Andrews

# General $\boldsymbol{\Sigma}$: rotated ellipses

$$\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{21} & \sigma_2^2 \end{bmatrix} \quad d_{\boldsymbol{\Sigma}}(\boldsymbol{x}; \mathbf{0}) \text{ are rotated ellipses}$$



$$\boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.4 \\ 0.4 & 1 \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & -0.4 \\ -0.4 & 1 \end{bmatrix} \qquad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & -0.9 \\ -0.9 & 1 \end{bmatrix}$$

# MLE of multivariate Gaussian

Given $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \ldots, \mathbf{x}^{(m)}\}$, assume $\mathbf{x}^{(i)} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$; the goal is to estimate

$$\boldsymbol{\mu}, \boldsymbol{\Sigma}$$

The log likelihood is:

$$\mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log P(\{\mathbf{x}^{(i)}\}_1^m | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^{m} \log N(\mathbf{x}^{(i)}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

The MLE is defined as usual:

$$\boldsymbol{\mu}_{ML}, \boldsymbol{\Sigma}_{ML} = \underset{\boldsymbol{\mu}, \boldsymbol{\Sigma}}{\operatorname{argmax}} \, \mathcal{L}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
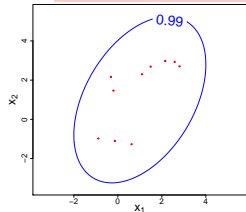
Take derivative and set to zero; after some tedious steps, the solution is:

$$\boldsymbol{\mu}_{ML} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{x}^{(i)}, \quad \boldsymbol{\Sigma}_{ML} = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{ML})(\mathbf{x}^{(i)} - \boldsymbol{\mu}_{ML})^T$$
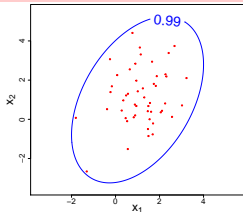
# Example: MLE of MV Gaussian

True parameters: $\boldsymbol{\mu} = [1,1]^T$, $\quad \boldsymbol{\Sigma} = \begin{bmatrix} 1 & 0.6 \\ 0.6 & 2 \end{bmatrix}$

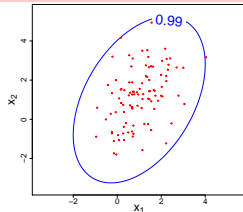- MLE is **consistent**: it converges to the truth with enough data



(sample size) : 10     (sample size) : 50     (sample size) : 100     (sample size) : 5000

$\boldsymbol{\mu}_{ML} = [0.92, 1.39]^T$   $\boldsymbol{\mu}_{ML} = [1.16, 1.16]^T$   $\boldsymbol{\mu}_{ML} = [1.04, 1.19]^T$   $\boldsymbol{\mu}_{ML} = [0.99, 0.99]^T$

$\boldsymbol{\Sigma}_{ML} = \begin{bmatrix} 1.54 & 1.45 \\ 1.45 & 2.86 \end{bmatrix}$   $\boldsymbol{\Sigma}_{ML} = \begin{bmatrix} 1.01 & 0.31 \\ 0.31 & 1.98 \end{bmatrix}$   $\boldsymbol{\Sigma}_{ML} = \begin{bmatrix} 0.87 & 0.55 \\ 0.55 & 2.01 \end{bmatrix}$   $\boldsymbol{\Sigma}_{ML} = \begin{bmatrix} 0.98 & 0.6 \\ 0.6 & 2.05 \end{bmatrix}$

# Finite mixture model

# EM for mixture of Gaussians

# Revisit K-means

University of
St Andrews

# Demonstration

University of
St Andrews

# How to decide $K$

# EM for general mixture

# EM as a general algorithm

# Review: expectation

**Expection** of a r.v. is defined as

$$\mathbb{E}[g(X)] = \sum_x g(x)P(x) \text{ or } \mathbb{E}[g(X)] = \int g(x)P(x)dx$$

- $\mathbb{E}[a] = a$ ($a$ is a constant)
- linearity: $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$
- $\mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X]$: as $\mathbb{E}[X]$ is a constant (the randomness has been integrated out)

Interpretation of Expectation: sample mean of a very large sample

$$\mathbb{E}[g(X)] = \frac{1}{m}\sum_{i=1}^{m} g(x^{(i)}); \quad m \to \infty$$

- limit of the sample average of $\{x^{(1)}, \ldots, x^{(m)}\}$ and $x^{(i)} \sim P(X)$

# Review: varaiance covariance

**Variance** of a r.v. is defined as

$$\mathrm{Var}[g(X)] = \mathbb{E}[(g(X) - \mathbb{E}[g(X)])^2] = \mathbb{E}[g(X)^2] - \mathbb{E}[g(X)]^2$$

- $\mathrm{Var}[aX] = a^2\mathrm{Var}[X]$
- measures the spread of the distribution around the mean $\mathbb{E}[g(X)]$

# Example

$X$ is a Bernoulli r.v. with parameter $p = 0.5$; what is $\mathbb{E}[X]$?
- $\mathbb{E}[X] = 1 \times P(X = 1) + 0 \times P(X = 0) = p = 0.5$;

$Y$ is a Binomial r.v. with $N = 10, p = 0.5$, what is $\mathbb{E}[Y]$?
- $Y = \sum_{i=1}^{N} X = N \times X$
- $\mathbb{E}[Y] = \mathbb{E}[N \times X] = N \times \mathbb{E}[X] = N \times p = 5$
- interpretation: you expect to see 5 successes out of 10 (on average the result is 5 if you repeat the experiment a lot of times)