

Introduction to Machine Learning
Summer 2018
Final Exam Solution

1 ML

$$\begin{aligned}\ell(\mu, \sigma^2) &= \log \left(\prod_{i=1}^N p_X(x_i; \mu, \sigma^2) \right) \\ &= \sum_{i=1}^N \left(\log \left(\frac{1}{x_i} \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2\sigma^2} (\log(x_i) - \mu)^2 \right) \quad x_i > 0\end{aligned}$$

$$\hat{\mu} = \arg \max_{\mu} \ell(\mu, \sigma^2)$$

$$\frac{\partial}{\partial \mu} \ell(\mu, \sigma^2) = 0$$

$$\sum_{i=1}^N (\log(x_i) - \mu) = 0$$

$$\mu = \frac{1}{N} \sum_{i=1}^N \log(x_i)$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum_{i=1}^N \log(x_i)$$

$$\Rightarrow \exp(\hat{\mu}) = \exp \left(\frac{1}{N} \sum_{i=1}^N \log(x_i) \right) = \exp \left(\log \left(\left(\prod_{i=1}^N x_i \right)^{\frac{1}{N}} \right) \right) = \left(\prod_{i=1}^N x_i \right)^{\frac{1}{N}}$$

We obtain the geometric mean.

For $\{x_1 = 2, x_2 = 3\}$ we have:

$$\exp(\hat{\mu}) = \sqrt{2 \cdot 3} = \sqrt{6}$$

2 MAP

$$\begin{aligned}
\hat{\lambda}_{MAP} &= \arg \max_{\lambda} p(\lambda | \{x_i\}) \\
&= \arg \max_{\lambda} p(\{x_i\} | \lambda) p_{\lambda}(\lambda) \\
&= \arg \max_{\lambda} \prod_{i=1}^N \left(\frac{\lambda^{x_i}}{x_i!} e^{-\lambda} \right) \cdot \frac{\beta^{\alpha}}{\Gamma(\alpha)} \lambda^{\alpha-1} e^{-\beta\lambda} \\
&= \arg \max_{\lambda} \prod_{i=1}^N (\lambda^{x_i} e^{-\lambda}) \cdot \lambda^{\alpha-1} e^{-\beta\lambda} \\
&= \arg \max_{\lambda} \lambda^{\sum_{i=1}^N x_i + \alpha - 1} e^{-(N+\beta)\lambda} \\
&= \arg \max_{\lambda} \underbrace{\left(\sum_{i=1}^N x_i + \alpha - 1 \right) \log(\lambda) - (N + \beta) \lambda}_{\triangleq f(\lambda)}
\end{aligned}$$

$$f'(\lambda) = 0$$

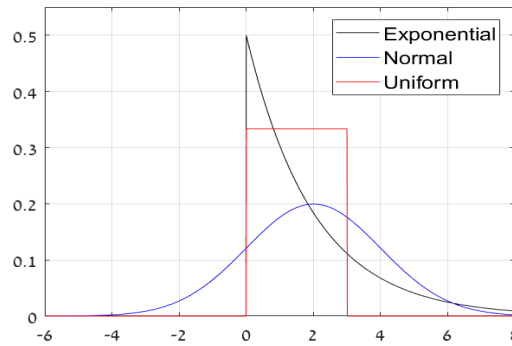
$$\left(\sum_{i=1}^N x_i + \alpha - 1 \right) \frac{1}{\lambda} - N - \beta = 0$$

$$\left(\sum_{i=1}^N x_i + \alpha - 1 \right) \frac{1}{\lambda} = N + \beta$$

$$\Rightarrow \hat{\lambda}_{MAP} = \frac{\sum_{i=1}^N x_i + \alpha - 1}{N + \beta} = \frac{N\bar{x} + \alpha - 1}{N + \beta} \xrightarrow{N \rightarrow \infty} \bar{x}$$

3 Bayes Classifier

The conditional probabilities densities are as follows



Consider the following cases:

- $x < 0$: In this case, only the conditional normal distribution has a positive density (thus, the prior distributions are irrelevant).
- $0 \leq x < 3$: The class u and g have the same prior probability and the conditional uniform density is larger than the conditional normal density $\left(\frac{1}{\sqrt{8\pi}} < \frac{1}{3} \right)$. Hence, we need compare between the uniform and the exponential posterior

probabilities

$$\begin{aligned}
0.4 \cdot \frac{1}{2} e^{-\frac{x}{2}} &\geq 0.3 \cdot \frac{1}{3} \\
\Rightarrow e^{-\frac{x}{2}} &\geq \frac{1}{2} \\
\Rightarrow -\frac{x}{2} &\geq -\ln(2) \\
\Rightarrow x &\leq 2 \ln(2) \approx 1.38.
\end{aligned}$$

- $x \geq 3$: In this case, the conditional uniform distribution is zero (its prior is irrelevant), hence, we compare the normal and exponential posterior probabilities:

$$\begin{aligned}
0.4 \cdot \frac{1}{2} e^{-\frac{x}{2}} &\geq 0.3 \cdot \frac{1}{\sqrt{8\pi}} e^{-\frac{(x-2)^2}{8}} \\
\Rightarrow e^{-\frac{x}{2}} &\geq \frac{1.5}{\sqrt{8\pi}} e^{-\frac{(x-2)^2}{8}} \\
\Rightarrow -\frac{x}{2} &\geq \ln(1.5) - \frac{1}{2} \ln(8\pi) - \frac{(x-2)^2}{8} \\
\Rightarrow x^2 - 4x + 4 - 4x - 8 \left(\ln(1.5) - \frac{1}{2} \ln(8\pi) \right) &\geq 0 \\
\Rightarrow x^2 - 8x + 4 - 8 \left(\ln(1.5) - \frac{1}{2} \ln(8\pi) \right) &\geq 0 \\
\Rightarrow x &\geq 5.532.
\end{aligned}$$

Note that the second solution $x \leq 2.468$ is irrelevant since we assumed $x \geq 3$.

Thus, the Bayes optimal classifier is given by

$$\hat{\omega} = \begin{cases} g & x < 0 \\ e & 0 \leq x < 1.38 \\ u & 1.38 \leq x < 3 \\ g & 3 \leq x < 5.532 \\ e & x \geq 5.532. \end{cases}$$

4 Histogram

$$\begin{aligned}
\mathbb{E}[\hat{p}_X(x_0)] &= \mathbb{E} \left[\frac{1}{N} \cdot \frac{1}{|R_k|} \sum_{i=1}^N \mathbf{I}\{x_i \in R_k\} \right] \\
&= \frac{1}{N} \cdot \frac{1}{b-a} \sum_{i=1}^N \mathbb{E}[\mathbf{I}\{x_i \in R_k\}] \\
&= \frac{1}{N} \cdot \frac{1}{b-a} \sum_{i=1}^N \Pr\{x_i \in R_k\} \\
&= \frac{1}{b-a} \Pr\{a < x_1 \leq b\} \\
&= \frac{F_X(b) - F_X(a)}{b-a}
\end{aligned}$$

5 PCA I

1.

$$\begin{aligned}
\|\mathbf{y}_i - \mathbf{y}_j\|_2 &= \left\| \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) - \mathbf{U}^T (\mathbf{x}_j - \boldsymbol{\mu}_x) \right\|_2 \\
&= \left\| \mathbf{U}^T \mathbf{x}_i - \mathbf{U}^T \mathbf{x}_j \right\|_2 \\
&= \left\| \mathbf{U}^T (\mathbf{x}_i - \mathbf{x}_j) \right\|_2 \\
&= \|\mathbf{x}_i - \mathbf{x}_j\|_2
\end{aligned}$$

2.

$$\boldsymbol{\mu}_y = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i = \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) = \mathbf{U}^T \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) = 0$$

$$\begin{aligned}
\boldsymbol{\Sigma}_{yy} &= \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i \mathbf{y}_i^T \\
&= \frac{1}{N} \sum_{i=1}^N \mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) \left(\mathbf{U}^T (\mathbf{x}_i - \boldsymbol{\mu}_x) \right)^T \\
&= \mathbf{U}^T \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_x) (\mathbf{x}_i - \boldsymbol{\mu}_x)^T \mathbf{U} \\
&= \mathbf{U}^T \boldsymbol{\Sigma}_x \mathbf{U} \\
&= \mathbf{U}^T \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{U} \\
&= \boldsymbol{\Lambda}
\end{aligned}$$

6 PCA II

As we saw in the tutorial, the new empirical covariance matrix is given

$$\begin{aligned}
\tilde{\boldsymbol{\Sigma}}_x &= \mathbf{V} \boldsymbol{\Sigma}_x \mathbf{V}^T \\
&= \mathbf{V} \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^T \mathbf{V}^T \\
&= (\mathbf{V} \mathbf{U}) \boldsymbol{\Lambda} (\mathbf{V} \mathbf{U})^T \\
&= \tilde{\mathbf{U}} \boldsymbol{\Lambda} \tilde{\mathbf{U}}^T
\end{aligned}$$

Therefore, the new principle components are given by $\tilde{\mathbf{U}} = \mathbf{V} \mathbf{U}$. Denote by $\tilde{\mathbf{U}}_m$ the first m principle components, we have that $\tilde{\mathbf{U}}_m = \mathbf{V} \mathbf{U}_m$, hence

$$\begin{aligned}
\tilde{\mathbf{y}}_i &= \tilde{\mathbf{U}}_m^T \tilde{\mathbf{x}}_i \\
&= \mathbf{U}_m^T \underbrace{\mathbf{V}^T \mathbf{V}}_{=I} \mathbf{x}_i \\
&= \mathbf{U}_m^T \mathbf{x}_i \\
&= \mathbf{y}_i
\end{aligned}$$

Thus, we can conclude the applying an orthonormal linear transformation does not change the representations.

7 K-Means

$$J_0 = \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

$$J_1 = \sum_{k=1}^K \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{m}_k\|_2^2$$

where:

$$\mathbf{m}_k = \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_i \in \mathcal{C}_k} \mathbf{x}_i$$

For simplicity, we consider a single cluster:

$$J_0 = \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2$$

$$J_1 = \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{m}_k\|_2^2$$

$$\begin{aligned} J_0 &= \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \boldsymbol{\mu}_k\|_2^2 \\ &= \sum_{\mathbf{x}_i \in \mathcal{C}_k} \|\mathbf{x}_i - \mathbf{m}_k + \mathbf{m}_k - \boldsymbol{\mu}_k\|_2^2 \\ &= \sum_{\mathbf{x}_i \in \mathcal{C}_k} \left(\|\mathbf{x}_i - \mathbf{m}_k\|_2^2 + \|\mathbf{m}_k - \boldsymbol{\mu}_k\|_2^2 + 2(\mathbf{x}_i - \mathbf{m}_k)^T (\mathbf{m}_k - \boldsymbol{\mu}_k) \right) \\ &= J_1 + 2(\mathbf{m}_k - \boldsymbol{\mu}_k)^T \underbrace{\sum_{\mathbf{x}_i \in \mathcal{C}_k} (\mathbf{x}_i - \mathbf{m}_k)}_{=0} + |\mathcal{C}_k| \|\mathbf{m}_k - \boldsymbol{\mu}_k\|_2^2 \\ &= J_1 + \underbrace{|\mathcal{C}_k| \|\mathbf{m}_k - \boldsymbol{\mu}_k\|_2^2}_{\geq 0} \\ &\geq J_1 \end{aligned}$$

$$\Rightarrow J_1 \leq J_0$$

8 Perceptron

8.1

The lower bound is zero (or one) iterations if the initial guess is already providing perfect classification.
For example:

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

and:

$$\mathbf{w}_1 = \begin{bmatrix} 10 \\ 0 \end{bmatrix}$$

$$\Rightarrow \text{sign}(\mathbf{w}_1^T \mathbf{x}_1) = \text{sign}(10) = 1 = y_1$$

8.2

The upper bound is 11 (or 12 if $\text{sign}(0) \neq 1$) iterations.

For example:

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

and

$$\mathbf{w}_1 = \begin{bmatrix} -10 \\ 0 \end{bmatrix}$$

$$\Rightarrow \text{sign}(\mathbf{w}_1^T \mathbf{x}_1) = \text{sign}(-10) = -1$$

So after one iteration we have:

$$\mathbf{w}_2 = \mathbf{w}_1 + \mathbf{x}_1 = \begin{bmatrix} -10 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -9 \\ 0 \end{bmatrix}$$

$$\Rightarrow \mathbf{w}_3 = \begin{bmatrix} -8 \\ 0 \end{bmatrix}$$

\vdots

$$\mathbf{w}_{11} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\mathbf{w}_{12} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\Rightarrow \text{sign}(\mathbf{w}_{11}^T \mathbf{x}_1) = \text{sign}(1) = 1 = y_1$$

9 Regression

1.

$$\begin{aligned} dL &= (-X dw)^T A(y - Xw) - (y - Xw)^T AX dw \\ &= -(y - Xw)^T (A^T + A) X dw \\ &= \left\langle -X^T (A^T + A) (y - Xw), dw \right\rangle \\ &\Rightarrow \nabla_w L = -X^T (A^T + A) (y - Xw) \end{aligned}$$

2.

$$\mathbf{w}_{k+1} = \mathbf{w}_0 - \mu \nabla_w L = \mathbf{w}_0 + \mu X^T (A^T + A) (y - Xw_0)$$

10 Kernel function

$$\begin{aligned} k(x, z) &= x^T A z \\ &= x^T U \Lambda U^T z \\ &= x^T U \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} U^T z \\ &= \left\langle \Lambda^{\frac{1}{2}} U^T x, \Lambda^{\frac{1}{2}} U^T z \right\rangle \end{aligned}$$

$$\Rightarrow \phi(x) = \Lambda^{\frac{1}{2}} U^T x$$