



Machine Learning 1

Winter semester 2016/17

Group APXNLE

Exercise 2

Members

Jing Li[387272] jing.li.1@campus.tu-berlin.de

Kumar Awanish[386697] k.awanish@campus.tu-berlin.de

Manjiao Xu[386498] manjiao.xu@campus.tu-berlin.de

Rudresha Gulaganjihalli Parameshappa[386642]

Gulaganjihalliparameshappa@campus.tu-berlin.de

Sonali Nayak[386995] sonali.nayak@campus.tu-berlin.de

Maximilian Ernst[364862] maximilian.ernst@campus.tu-berlin.de

Exercise 1: Maximum-Likelihood Estimation

1. (a). $p(x) = \int_{-\infty}^{+\infty} p(x, y) dy$ (supported on \mathbb{R}_+^2)

$$\begin{aligned}
 &= \int_0^{+\infty} \lambda \eta e^{-\lambda x - \eta y} dy \\
 &= \lambda \eta \int_0^{+\infty} e^{-\lambda x} * e^{-\eta y} dy \\
 &= \lambda \eta e^{-\lambda x} \int_0^{+\infty} e^{-\eta y} dy \\
 &= \lambda \eta e^{-\lambda x} * \left(-\frac{1}{\eta} e^{-\eta y} \Big|_0^{+\infty} \right) \\
 &= \lambda \eta e^{-\lambda x} * \left(-\frac{1}{\eta} (0 - 1) \right) \\
 &= \lambda \eta e^{-\lambda x} * \frac{1}{\eta} \\
 &= \lambda e^{-\lambda x}
 \end{aligned}$$

$$\begin{aligned}
 p(y) &= \int_{-\infty}^{+\infty} p(x, y) dx \\
 &= \int_0^{+\infty} \lambda \eta e^{-\lambda x - \eta y} dx \\
 &= \lambda \eta \int_0^{+\infty} e^{-\lambda x} * e^{-\eta y} dx \\
 &= \lambda \eta e^{-\eta y} \int_0^{+\infty} e^{-\lambda x} dx \\
 &= \lambda \eta e^{-\eta y} * \left(-\frac{1}{\lambda} e^{-\lambda x} \Big|_0^{+\infty} \right) \\
 &= \lambda \eta e^{-\eta y} * \left(-\frac{1}{\lambda} (0 - 1) \right) \\
 &= \lambda \eta e^{-\eta y} * 1 \\
 &= \eta e^{-\eta y}
 \end{aligned}$$

$$\begin{aligned}
 \therefore p(x) * p(y) &= \lambda e^{-\lambda x} * \eta e^{-\eta y} \\
 &= \lambda \eta e^{-\lambda x - \eta y} \\
 &= p(x, y)
 \end{aligned}$$

$\therefore x$ and y are independent.

(b) $L(\lambda) = \prod_{i=1}^n p(x_i, y_i)$

$$\begin{aligned}
 &= \prod_{i=1}^n \lambda \eta e^{-\lambda x_i - \eta y_i} \\
 &= \lambda^n \eta^n e^{-\lambda \sum_{i=1}^n x_i - \eta \sum_{i=1}^n y_i}
 \end{aligned}$$

$$\ln L(\lambda) = n \ln \lambda + n \ln \eta - \lambda \sum_{i=1}^n x_i - \eta \sum_{i=1}^n y_i$$

$$\begin{aligned}
 \frac{d \ln L(\lambda)}{d \lambda} &= \frac{n}{\lambda} + 0 - \sum_{i=1}^n x_i - 0 \\
 &= \frac{n}{\lambda} - \sum_{i=1}^n x_i
 \end{aligned}$$

$$\frac{n}{\lambda} - \sum_{i=1}^n x_i = 0 \quad \text{according to } \sum_{i=1}^n x_i = n \bar{x}$$

$$\begin{aligned}
 \frac{n}{\lambda} - n \bar{x} &= 0 \\
 \hat{\lambda} &= \frac{1}{\bar{x}}
 \end{aligned}$$

(c) $\eta = \frac{1}{\lambda}$

$$\begin{aligned}
 L(\lambda) &= \prod_{i=1}^n p(x_i, y_i) \\
 &= \prod_{i=1}^n \lambda * \frac{1}{\lambda} * e^{-\lambda x_i - \frac{1}{\lambda} y_i} \\
 &= e^{-\lambda \sum_{i=1}^n x_i - \frac{1}{\lambda} \sum_{i=1}^n y_i}
 \end{aligned}$$

$$\ln L(\lambda) = -\lambda \sum_{i=1}^n x_i - \frac{1}{\lambda} \sum_{i=1}^n y_i$$

$$\begin{aligned}
 \frac{d \ln L(\lambda)}{d \lambda} &= -\sum_{i=1}^n x_i - \left(-\frac{1}{\lambda^2} \sum_{i=1}^n y_i \right) \\
 &= -\sum_{i=1}^n x_i + \frac{1}{\lambda^2} \sum_{i=1}^n y_i
 \end{aligned}$$

$$\frac{1}{\lambda^2} \sum_{i=1}^n y_i - \sum_{i=1}^n x_i = 0$$

$$\frac{1}{\lambda^2} \bar{y} - n \bar{x} = 0$$

$$\frac{\bar{y}}{\lambda^2} = \bar{x}$$

$$\lambda^2 = \bar{y} / \bar{x}$$

$$\hat{\lambda} = \pm \sqrt{\bar{y} / \bar{x}}$$

we have $\lambda > 0$

$$\text{so } \hat{\lambda} = \sqrt{\bar{y} / \bar{x}}$$

$$(d). \eta = 1 - \lambda$$

$$\begin{aligned} L(\lambda) &= \prod_{i=1}^n p(x_i, y_i) \\ &= \prod_{i=1}^n \lambda (1-\lambda) e^{-\lambda x_i - (1-\lambda)y_i} \\ &= \lambda^n (1-\lambda)^n e^{-\lambda \sum_{i=1}^n x_i - (1-\lambda) \sum_{i=1}^n y_i} \end{aligned}$$

$$\begin{aligned} \ln L(\lambda) &= n \ln \lambda + n \ln (1-\lambda) - \lambda \sum_{i=1}^n x_i - (1-\lambda) \sum_{i=1}^n y_i \\ &= n \ln \lambda + n \ln (1-\lambda) - \lambda \sum_{i=1}^n x_i - (1-\lambda) \sum_{i=1}^n y_i \end{aligned}$$

$$\frac{d \ln L(\lambda)}{d \lambda} = \frac{n}{\lambda} + (-\frac{n}{1-\lambda}) - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$$

$$\frac{n}{\lambda} - \frac{n}{1-\lambda} - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i = 0$$

$$\frac{n}{\hat{\lambda}} - \frac{n}{1-\hat{\lambda}} - n\bar{x} + n\bar{y} = 0$$

$$1-\hat{\lambda} - \hat{\lambda} - \hat{\lambda}(1-\hat{\lambda})\bar{x} + \hat{\lambda}(1-\hat{\lambda})\bar{y} = 0$$

$$(\bar{y} - \bar{x})\hat{\lambda}^2 + (\bar{x} - \bar{y} + 2)\hat{\lambda} - 1 = 0$$

$$\hat{\lambda} = \frac{-(\bar{x} - \bar{y} + 2) \pm \sqrt{(\bar{x} - \bar{y} + 2)^2 - 4(\bar{y} - \bar{x})(-1)}}{2(\bar{y} - \bar{x})}$$

$$= \frac{\bar{y} - \bar{x} - 2 \pm \sqrt{(\bar{x} - \bar{y} + 2)^2 + 4(\bar{y} - \bar{x})}}{2(\bar{y} - \bar{x})}$$

$$\text{we have } \lambda > 0 \quad \eta > 0$$

$$\eta = 1 - \lambda > 0 \quad \lambda < 1$$

$$\text{so } 0 < \lambda < 1$$

Exercise 2: Multiple Linear Regression

Given: Multiple linear regression problem $y = x^T \beta + \epsilon$, with:

$x \in \mathbb{R}^d$: the predictor variables,

$y \in \mathbb{R}$: the response variable, and

$\beta \in \mathbb{R}^d$: the linear regression coefficients.

Data: $\mathcal{D} = ((x_1, y_1), \dots, (x_N, y_N))$ of N independent data pairs.

Note $Y := (y_1, \dots, y_N)^T \in \mathbb{R}^N$, $\tilde{\epsilon} = (\epsilon_1, \dots, \epsilon_N)^T \in \mathbb{R}^N$ and

$X = (x_1, \dots, x_N)^T \in \mathbb{R}^{N \times d}$.

(a) Show that $\hat{\beta} := (X^T X)^{-1} X^T Y \sim \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1})$

Proof. Since $y = x^T \beta + \epsilon$, there is

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} x_1^1 \dots x_d^1 \\ \vdots \\ x_1^N \dots x_d^N \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{pmatrix} = X\beta + \tilde{\epsilon}.$$

We know that $\epsilon \sim \mathcal{N}(0, \sigma^2)$, thus $\tilde{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_N \end{pmatrix}$ is multivariate Gaussian

distributed and $\tilde{\epsilon} \sim \mathcal{N}(0, \sigma^2 I)$ holds, where I is the identity matrix.

It follows that with the **affine transformation**, $Y = X\beta + \tilde{\epsilon}$ is also multivariate Gaussian distributed: $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$. This is well-defined since I is not singular.

Because of the linear property of the expected value, it follows that

$$\begin{aligned} \mathbb{E}[\hat{\beta}] &= \mathbb{E}[(X^T X)^{-1} X^T Y] \\ &= (X^T X)^{-1} X^T \mathbb{E}[Y] \\ &= (X^T X)^{-1} X^T \mathbb{E}[X\beta + \tilde{\epsilon}] \\ &= (X^T X)^{-1} (X^T X) \beta + (X^T X)^{-1} X^T \mathbb{E}[\tilde{\epsilon}] \\ &= \beta + 0 \\ &= \beta \end{aligned}$$

and

$$\begin{aligned}
\Sigma[\hat{\beta}] &= \Sigma(X^T X)^{-1} X^T Y \\
&= ((X^T X)^{-1} X^T) \cdot \Sigma[Y] \cdot ((X^T X)^{-1} X^T)^T \\
&= ((X^T X)^{-1} X^T) \cdot \sigma^2 I \cdot X (X^T X)^{-T} \\
&= ((X^T X)^{-1} X^T) \cdot \sigma^2 \cdot X (X^T X)^{-T} \\
&= \sigma^2 (X^T X)^{-1} X^T X (X^T X)^{-T} \\
&= \sigma^2 (X^T X)^{-T} \\
&= \sigma^2 (X^T X)^{-1}
\end{aligned}$$

As the affine transformed of the multivariate Gaussian distributed random variable Y , $\hat{\beta}$ is also Gaussian distributed. It has mean β and covariance matrix $\sigma^2 (X^T X)^{-1}$. The covariance is well defined since $(X^T X)^{-1}$ is not singular. □

(b) Knowing the full distribution of $\hat{\beta}$ makes it possible for us to choose the variable β from the whole real-axis \mathbb{R} , i.e. the estimated value is no longer one single value anymore, but rather a set of real numbers or an interval.

(c) Given new data point x_* , we want to predict the response for x_* . We know that the response is $y_* := x_*^T \hat{\beta} = x_*^T (X^T X)^{-1} X^T Y$, where $\hat{\beta} \sim \mathcal{N}(\beta, \sigma^2 (X^T X)^{-1})$.

Thus, as the result of an affine transformation, y_* is 1-dim normally distributed and has expectation $x_*^T \beta$ and variance $\sigma^2 x_*^T (X^T X)^{-1} x_*$.

(d) Now that we have the distribution of the new predicted response, we are able to see with which values the response is inclined to appear. We can also deduce other information from it, e.g. constructing the confidence interval of the predicted response.