

CPSC340



Probabilistic linear prediction



Nando de Freitas October, 2012 University of British Columbia

Outline of the lecture

In this lecture, we formulate the problem of linear prediction using probabilities. In doing so, we find out that the maximum likelihood estimate coincides with the least squares estimate. The goal of the lecture is for you to learn:

Multivariate Gaussian distributions
 How to formulate the likelihood for linear regression
 Computing the maximum likelihood estimates for linear regression.

Univariate Gaussian distribution

The probability density function (pdf) of a Gaussian distribution is given by

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}.$$

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

where μ is the mean or center of mass and σ^2 is the variance.



Multivariate Gaussian distribution

Let $\mathbf{y} \in \mathbb{R}^{n \times 1}$, then pdf of an n-dimensional Gaussian is given by

 $p(\mathbf{y}) = |2\pi \mathbf{\Sigma}|^{-1/2} e^{-\frac{1}{2}(\mathbf{y}-\boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{y}-\boldsymbol{\mu})}, \quad \mathbf{\hat{z}}$ N=2 where $\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} \mathbb{E}(y_1) \\ \vdots \\ \mathbb{E}(y_n) \end{pmatrix}$ 2=21 and $\mathbf{S}_{\mathbf{u}} = \mathbf{G}_{\mathbf{u}} \qquad \mathbf{\Sigma} = \begin{pmatrix} \sigma_{11} \cdots \sigma_{1n} \\ \cdots \\ \sigma_{n1} \cdots \sigma_{nn} \end{pmatrix} = \mathbb{E}[(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})^T]$ $\sum_{i=1}^{n} \begin{pmatrix} \overline{b}_{11} & \overline{b}_{12} \\ \overline{b}_{12} & \overline{b}_{22} \end{pmatrix} \xrightarrow{A}$



Let us assume that we have n=3 data points $y_1 = 1$, $y_2 = 0.5$, $y_3 = 1.5$, which we assume independent and Gaussian distributed with unknown mean θ and variance 1. That is,

 $y_i = \mathcal{N}(\theta, 1) = \theta + \mathcal{N}(\theta, 1)$

with likelihood $P(y_1y_2y_3|\theta) = P(y_1|\theta) P(y_1|\theta) P(y_3|\theta)$. Consider the following two cases. Clearly the one on the left has higher likelihood (higher green bars). Finding the θ that maximizes the likelihood is equivalent to moving the Gaussian to the left and right until the product of 3 green bars (likelihood) is maximized.

The likelihood for linear regression

Let us assume that each label y_i is Gaussian distributed with mean $x_i^T \theta$ and variance σ^2 , which in short we write as:

$$\mathbf{y}_{i} = \mathcal{N}(\mathbf{x}_{i}^{T}\boldsymbol{\theta}, \sigma^{2}) = \mathbf{x}_{i}^{T}\boldsymbol{\theta} + \mathcal{N}(\boldsymbol{\theta}, \sigma^{2})$$

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma) = \prod_{i=1}^{n} p(y_{i}|\mathbf{x}_{i}, \boldsymbol{\theta}, \sigma).$$

$$= \prod_{i=1}^{n} (2\pi\sigma^{2})^{-1/2} e^{-\frac{1}{2\sigma^{2}}(y_{i} - \mathbf{x}_{i}^{T}\boldsymbol{\theta})^{2}}$$

$$= (2\pi\sigma^{2})^{-n/2} e^{-\frac{1}{2\sigma^{2}}\sum_{i=1}^{n}(y_{i} - \mathbf{x}_{i}^{T}\boldsymbol{\theta})^{2}}$$

$$= (2\pi\sigma^{2})^{-n/2} e^{-\frac{1}{2\sigma^{2}}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^{T}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}$$

Maximum likelihood

The maximum likelihood estimate (MLE) of $\boldsymbol{\theta}$ is obtained by taking the derivative of the log-likelihood, $\log p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma)$. The goal is to maximize the likelihood of seeing the training data \mathbf{y} by modifying the parameters $(\boldsymbol{\theta}, \sigma)$.

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}, \sigma) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}$$
$$\ell(\mathbf{y}, \mathbf{y}) = \log P(\mathbf{y}|\mathbf{X}, \mathbf{\theta}, \mathbf{g})$$
$$= -\frac{n}{2} \log (2\pi) - \frac{n}{2} \log^{2} - \frac{1}{2g^2} (\mathbf{y} - \mathbf{X}\mathbf{\theta})^T (\mathbf{y} - \mathbf{X}\mathbf{\theta})$$

The ML estimate of θ is: (assume G known) $\frac{\partial}{\partial \Theta} \ell(\Theta, G^2) = -\frac{1}{2G^2} \frac{\partial}{\partial G} \left[\gamma - \chi \Theta \right] \left[\gamma - \chi \Theta \right]$ $= -\frac{1}{2G^2} \frac{\partial}{\partial G} \left(\gamma^T \gamma - 2 \gamma^T \times \Theta + 2 \chi^T \times \Theta \right)$ (just as before) $= -\frac{1}{c^{2}} \left(-\gamma^{T} X \Theta + \chi^{T} X \Theta \right)$ Equating to zero, yields $\chi^{\mathsf{T}} \chi \ominus = \chi^{\mathsf{T}} \chi$ $\hat{\boldsymbol{\Theta}}_{MI} = (\boldsymbol{X}^{\mathsf{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{Y}$ (same as least Squares) The ML estimate of σ is: (assume \odot known)

$$\frac{\partial}{\partial G} \ell(\Theta, G^{2}) = \frac{\partial}{\partial G} \left(-\frac{n}{2} \log G^{2} - \frac{1}{2G^{2}} \left(Y - X \Theta \right)^{T} \left(Y - X \Theta \right) \right)$$

$$= \left(-\frac{n}{2}\right)\frac{2\varsigma}{c^2} - \frac{1}{2}\left(-2\right)\frac{1}{c^3}\left(\gamma - \times \Theta\right)^{T}\left(\gamma - \times \Theta\right)$$

$$= -\frac{n}{G} + \frac{1}{G^3} (Y - X\Theta)^T (Y - X\Theta)$$

Equating to zero, yields

$$G^2 = \frac{1}{N} (Y - X\Theta)^T (Y - X\Theta) = \frac{1}{N} \sum_{i=1}^{N} (Y_i - X_i^T\Theta)^2$$

us expected, it is the estimator of variance.

Making predictions

The ML plugin prediction, given the training data D=(X, y), for a new input x_* and known σ^2 is given by:

The DAG for linear regression

Frequentist model

 $P(Y|X,G^2,\Theta)P(\Theta)P(G^2)$

Next lecture

In the next lecture, we introduce ridge regression and the Bayesian learning approach for linear predictive models.