

Theory and Concepts of Regression

Abstract :

- × We make the connection between linear regression and a Gaussian distribution
- × We show that the MSE is a maximum likelihood estimation procedure

Source :

- Kevin P. Murphy: Machine Learning: A Probabilistic Perspective
- Runkler, Skiena

Properties of Maximum Likelihood:

Consider $\mathbf{X} = \{\vec{x}_1, \dots, \vec{x}_N\}$ N examples drawn from $P_{\text{data}}(\vec{x})$. Let $p_{\text{MODEL}}(\vec{x}, \vec{\theta})$ be a parametric model used to estimate $P_{\text{data}}(\vec{x})$.

The maximum likelihood estimator for $\vec{\theta}$ is then defined as

$$\begin{aligned}\vec{\theta}_{\text{ML}} &= \arg \max_{\vec{\theta}} [p_{\text{MODEL}}(\mathbf{X}; \vec{\theta})] \\ &= \arg \max_{\vec{\theta}} \left[\prod_{i=1}^N p_{\text{MODEL}}(\vec{x}_i; \vec{\theta}) \right]\end{aligned}$$

For convenience take the log (does not change its arg max):

$$\vec{\theta}_{\text{ML}} = \arg \max_{\vec{\theta}} \left[\sum_{i=1}^N \log p_{\text{MODEL}}(\vec{x}_i; \vec{\theta}) \right]$$

Interpretation: Divide the formula by N to arrive at a expectation value:

$$\vec{\theta}_{ML} = \arg \max_{\vec{\theta}} \mathbb{E} [\log p_{MODEL}(\vec{x}, \vec{\theta})]$$

One way to interpret it is to view it as minimizing the dissimilarity between the empirical distribution \hat{p}_{DATA} and the model distribution p_{MODEL} :

$$\mathbb{E} [\log \hat{p}_{DATA}(\vec{x}) - \log p_{MODEL}(\vec{x})]$$

minimize the cross-entropy between the distributions

The mean squared error (MSE) is the cross-entropy between empirical data and a Gaussian model!

Asymptotically ($N \rightarrow \infty$) the MLE is the best estimator: ~~when~~ If the number of training examples would approach infinity, the MLE of a parameter converges to the true value of the parameter.

1.1 Linear regression

Linear regression formula:

$$\begin{aligned} y(\vec{x}) &= \vec{\omega}^T \cdot \vec{x} + \varepsilon = \\ &= \sum_{j=1}^D \omega_j x_j + \varepsilon \end{aligned}$$

with $\vec{\omega}$ the weight vector,
 ε the residual error (between
our linear predictions and the
true response)

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

We often assume $\varepsilon \sim \mathcal{N}(\mu, \sigma^2)$ Gaussian distributed

To make the connection between linear regression and Gaussians more explicit, we rewrite the model as

$$p(y | \vec{x}, \vec{\theta}) = \mathcal{N}(y | \underbrace{\mu(\vec{x})}_{\text{mean of the Gaussian}}, \sigma^2(\vec{x}))$$

with $\mu = \vec{\omega}^T \cdot \vec{x}$, $\sigma^2(\vec{x}) = \sigma^2$,
 $\vec{\theta} = (\vec{\omega}, \sigma^2)$

For 1-dimensional input:

$$\mu(\vec{x}) = \underbrace{\omega_0}_{\text{intercept/bias term}} + \underbrace{\omega_1}_{\text{slope}} \cdot x = \vec{\omega}^T \cdot \vec{x}$$

Non-linear relationships by replacing \vec{x} with some non-linear function $\phi(\vec{x})$:

$$p(y | \vec{x}, \vec{\theta}) = \mathcal{N}(y | \underbrace{\vec{\omega}^T \cdot \phi(\vec{x})}_{\text{basis function}}, \underbrace{\sigma^2}_{\text{constant variance}})$$

Example: Polynomial regression

$$\phi(x) = [1, x, x^2, \dots, x^d]$$

1.2. Logistic regression

Two changes for binary classification:

a) Gaussian \longrightarrow Bernoulli distribution

$$p(y | \vec{x}, \vec{\omega}) = \text{Ber}(y | \mu(\vec{x}))$$

↑
tossing a coin

where $\mu(\vec{x}) = \mathbb{E}(y(\vec{x})) = p(y=1 | \vec{x})$

b) Pass $\vec{\omega}^T \vec{x}$ through a function that ensures $0 \leq \mu(\vec{x}) \leq 1$

$$\mu(\vec{x}) = \text{sigm}(\vec{\omega}^T \cdot \vec{x}),$$

where $\text{sigm} = \text{sigmoid}/\text{logistic}/\text{logit}$ is defined as

$$\text{sigm}(\eta) = \frac{1}{1 + e^{-\eta}} = \frac{e^{\eta}}{e^{\eta} + 1}$$

Putting these two steps together we get

$$p(y | \vec{x}, \vec{w}) = \text{Ber}(y | \text{sigm}(\vec{w}^T \vec{x}))$$

This is called logistic regression due to its similarity to linear regression.

If we choose a threshold 0.5, we can induce a decision rule (linear decision boundary)

$$\hat{y}(x) = 1 \Leftrightarrow p(y=1 | \vec{x}) > 0.5$$

2. Maximum Likelihood Estimation

As discussed

$$p(y | \vec{x}, \vec{\theta}) = \mathcal{N}(y | \underbrace{\vec{\omega}^T \cdot \vec{x}}_{\mu(\vec{x})}, \sigma^2)$$

To estimate the parameters $\vec{\theta}$ we define the MLE (= maximum likelihood estimation):

$$\hat{\vec{\theta}} = \arg \max_{\vec{\theta}} [\log p(y | \vec{\theta})]$$

We can rewrite the log-likelihood as follows (assuming that the training examples are independent and identically distributed)

$$l(\vec{\theta}) = \log p(y | \vec{\theta}) = \sum_{i=1}^N \log p(y_i | \vec{x}_i, \vec{\theta})$$

Instead of maximizing the log-likelihood, we can minimize the negative log likelihood or NLL: more convenient for numerical evaluation.

$$NLL(\vec{\theta}) = - \sum_{i=1}^N \log p(y_i | \vec{x}_i, \vec{\theta})$$

Now we insert the definition of the Gaussian into the log-likelihood $l(\vec{\theta})$:

$$l(\vec{\theta}) = \sum_{i=1}^N \log \left[\left(\frac{1}{2\pi\sigma^2} \right)^{1/2} \exp \left(-\frac{1}{2\sigma^2} (y_i - \vec{w}^T \vec{x}_i)^2 \right) \right]$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \vec{w}^T \vec{x}_i)^2 - \frac{N}{2} \log(2\pi\sigma^2)$$

RSS = residual sum
of squares

$$RSS(w) = \|\epsilon\|_2^2$$

(ℓ_2 -norm)

$$MSE = \frac{RSS}{N} \quad \text{mean squared error}$$

The MLE (maximum likelihood estimation) for $\vec{\omega}$ is the one that minimizes the RSS, so this method is known as least squares. \Rightarrow MSE is a maximum likelihood estimation procedure

Next we derive the maximum likelihood estimation (MLE):

$$\begin{aligned} \text{NLL}(\vec{\omega}) &= \frac{1}{2} (\vec{y} - \vec{X} \vec{\omega})^T (\vec{y} - \vec{X} \vec{\omega}) \\ &= \frac{1}{2} \vec{\omega}^T \vec{X}^T \vec{X} \vec{\omega} - \vec{\omega}^T (\vec{X}^T \vec{y}) \\ &\quad + \frac{1}{2} \vec{y}^T \vec{y} \end{aligned}$$

$$\text{where } \vec{X}^T \vec{X} = \sum_{i=1}^N \vec{X}_i \vec{X}_i^T = \sum_{i=1}^N \begin{pmatrix} X_{i,1}^2 & \dots & X_{i,1} X_{i,D} \\ \vdots & & \vdots \\ X_{i,D} X_{i,1} & \dots & X_{i,D}^2 \end{pmatrix}$$

sum of squares matrix

$$\text{and } \vec{X}^T \vec{y} = \sum_{i=1}^N \vec{X}_i y_i$$

The gradient $g(\vec{w})$ is given as

$$g(\vec{w}) = [\vec{X}^T \cdot \vec{X} \cdot \vec{w} - \vec{X}^T \cdot \vec{y}] =$$

$$= \sum \vec{x}_i (\vec{w} \cdot \vec{x}_i - y_i)$$

Equation to zero ($g(\vec{w}) = 0$) gives

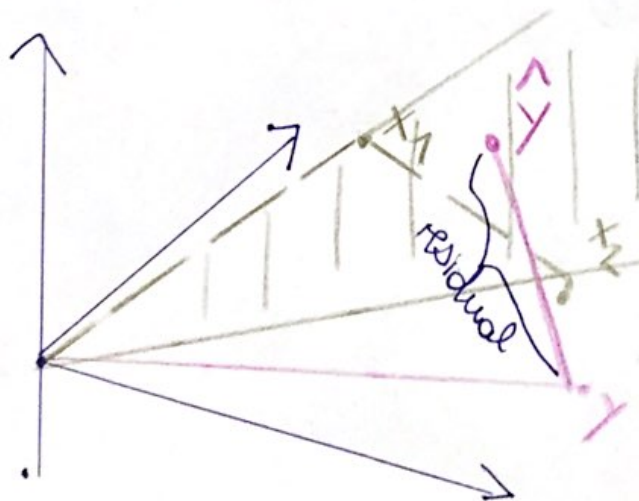
$$\boxed{\vec{X}^T \cdot \vec{X} \cdot \vec{w} = \vec{X}^T \cdot \vec{y}} \quad \text{normal equation}$$

The solution to this equation is

$$\boxed{\hat{w}_{OLS} = (\vec{X}^T \cdot \vec{X})^{-1} \cdot \vec{X}^T \cdot \vec{y}}$$

ordinary least squares solution

Graphical interpretation:



- x_1 and x_2 define a 2D plane
- y does not lie on that plane
- the orthogonal projection of y on that plane is \hat{y}

1-dim example:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \omega_0 - \omega_1 x_i)^2$$

$$\frac{\partial MSE}{\partial \omega_0} = -\frac{2}{N} \sum_{i=1}^N (y_i - \omega_0 - \omega_1 x_i) \stackrel{!}{=} 0$$

$$\Rightarrow \omega_0 = \bar{y} - \omega_1 \bar{x}$$

$$\Rightarrow MSE = \frac{1}{N} \sum (y_i - \bar{y} - \omega_1 (x_i - \bar{x}))^2$$

$$\frac{\partial MSE}{\partial \omega_1} = -\frac{2}{N} \sum (x_i - \bar{x})(y_i - \bar{y} - \omega_1 (x_i - \bar{x})) \stackrel{!}{=} 0$$

$$\Rightarrow \boxed{\omega_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{c_{ij}}{c_{jj}}}$$

Calculation example:

$$Y = \begin{bmatrix} 6 \\ 2 \\ 0 \\ 0 \\ 2 \end{bmatrix} \quad X = \begin{bmatrix} 4 & -2 \\ 1 & -1 \\ 0 & 0 \\ 1 & 1 \\ 4 & 2 \end{bmatrix}$$

$$\bar{y} = \frac{6+2+2}{5} = 2$$

⋮

$$\omega_0 = 2, \omega_1 = 1, \omega_2 = -1$$

3. Robust linear regression

To model the noise in regression models using a Gaussian distribution can be a bad choice when we have outliers in our fit. This is because squared error penalizes deviations quadratically, so points far from the line have more effect on the fit than points near the line.

for the response variable

Idea: Replace the Gaussian with a distribution that has heavy tails.

One possibility is to use the Laplace distrib:

$$p(y | \vec{x}, \vec{w}, b) = \text{Lap}(y | \vec{w} \cdot \vec{x}, b)$$

It assigns higher likelihood to outliers.

$$\sim \exp\left(-\frac{1}{b} |y - \vec{w} \cdot \vec{x}|\right)$$

$\left(\frac{1}{2b} e^{-\frac{|x-y|}{b}}\right)$

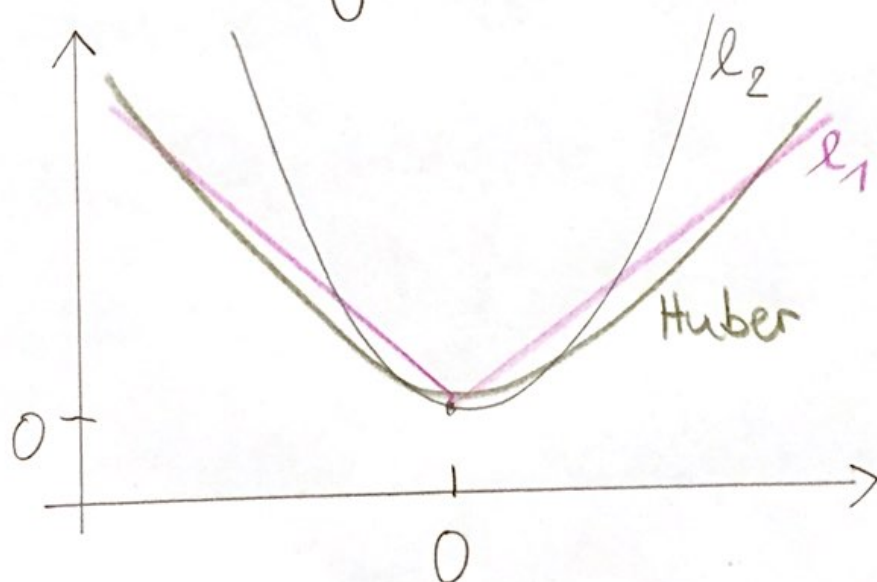
The robustness arises from using $|y - \vec{\omega}^T \cdot \vec{x}|$ instead of $(y - \vec{\omega}^T \cdot \vec{x})^2$.

$$\ell(\vec{\omega}) = \sum |y_i - \vec{\omega}^T \cdot \vec{x}_i| \quad \text{log-likelihood}$$

However, this non-linear objective function is hard to optimize. An alternative is to minimize the Huber

$$\text{loss} \quad L_H(r, \delta) = \begin{cases} r^2/2 & \text{if } |r| \leq \delta \\ \delta|r| - \frac{\delta^2}{2} & \text{if } |r| > \delta. \end{cases}$$

This is equivalent to ℓ_2 for error smaller than δ and ℓ_1 for larger errors. This loss function is everywhere differentiable.



4. Ridge Regression

Changing the log-likelihood $l(\vec{\theta})$ to:

$$\arg \max_{\vec{\omega}} \left\{ \sum_{i=1}^N \log \mathcal{N}(y_i | \omega_0 + \vec{\omega}^T \cdot \vec{x}_i, \sigma^2) + \sum_{j=1}^D \log \mathcal{N}(\omega_j | 0, \tau^2) \right\}$$

Zero-mean Gaussian prior
(to encourage parameters to be small)

$$\frac{P(A) \cdot P(B|A)}{P(B) \cdot P(A|B)}$$

This is equivalent to minimizing penalty

$$J(\vec{\omega}) = \frac{1}{N} \sum_{i=1}^N \underbrace{(y_i - (\omega_0 + \vec{\omega}^T \cdot \vec{x}_i))^2}_{\text{MSE}} + \lambda \|\vec{\omega}\|_2^2$$

with $\lambda = \frac{\sigma^2}{\tau^2}$ and $\|\vec{\omega}\|_2^2 = \vec{\omega}^T \vec{\omega} = \sum_{j=1}^D \omega_j^2$
squared two-norm.

The corresponding solution is given as

$$\hat{\omega}_{\text{Ridge}} = (\lambda \mathbb{1}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{penalized least squares}$$

Adding a Gaussian prior to the parameters of the model is called Ridge / l_2 regular or weight decay.

Note: $(\lambda \mathbb{1}_D + \mathbf{X}^T \mathbf{X})$ is easier to invert than $\mathbf{X}^T \mathbf{X}$, at least for large enough λ .

Note: In theory ridge regression is a biased MLE. A more "effective" approach is to increase the training data N . The approximation error should go to zero for $N \rightarrow \infty$.

5. Logistic Regression

$$p(y | \vec{x}, \vec{\omega}) = \text{Ber}(y | \text{sigm}(\vec{\omega}^T \cdot \vec{x}))$$

$$\begin{aligned} \text{NLL}(\vec{\omega}) &= - \sum_{i=1}^N \log[\mu_i^{\mathbb{1}(y_i=1)} \cdot (1-\mu_i)^{\mathbb{1}(y_i=0)}] \\ &= - \sum_{i=1}^N [y_i \log \mu_i + (1-y_i) \log(1-\mu_i)] \end{aligned}$$

Unlike for linear regression we can no longer write down the MLE in closed form. We need an optimization algorithm to compute it.

→ gradient descent