

Linear Regression



Regression problem

Given:

- observations

$$X = \{x_1, x_2, \dots, x_N\}, x_i \in \mathbb{R}^D$$

- targets

$$y = \{y_1, y_2, \dots, y_N\}, y_i \in \mathbb{R}$$

Find:

- mapping $f(\cdot)$ from inputs to targets

$$y_i \approx f(x_i)$$

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

↓
noise

$f(x) \rightarrow$ linear function

$$f_w(x_i) = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_D x_{iD}$$

$$= w_0 + w^T x_i$$

weight ↓ features of x_i

bias (offset)

→ absorb bias term

$$\tilde{x} = (1, x_1, \dots, x_D)^T$$

$$\tilde{w} = (w_0, w_1, \dots, w_D)^T$$

$$f_w(x) = \tilde{w}^T \tilde{x} \Rightarrow f_w(x) = w^T x$$

Loss function

Measures the "misfit" or error between our model (parametrized by w) and observed data $D = \{(x_i, y_i)\}_{i=1}^N$

- Least squares (LS)

$$E_{LS}(w) = \frac{1}{2} \sum_{i=1}^N (f_w(x_i) - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

Objective: find optimal weight vector w^* that minimizes the error

$$w^* = \arg \min_w E_{LS}(w) = \arg \min_w \frac{1}{2} \sum_{i=1}^N (x_i^T w - y_i)^2 = \arg \min_w \frac{1}{2} (x_w - y)^T (x_w - y)$$

To find the minimum of the loss $E(w)$, compute the gradient $\nabla_w E(w)$:

$$\begin{aligned}\nabla_w E_{LS}(w) &= \nabla_w \frac{1}{2} (Xw - y)^T (Xw - y) \\ &= \nabla_w \frac{1}{2} (w^T X^T X w - 2 w^T X^T y + y^T y) \\ &= X^T X w - X^T y\end{aligned}$$

Then set it to zero and solve for w to obtain the minimizer

$$X^T X w - X^T y \stackrel{!}{=} 0$$

$$\begin{aligned}w^* &= \underbrace{(X^T X)^{-1}}_{= X^+} X^T y \\ &= X^+ (Moore-Penrose pseudo inverse of X)\end{aligned}$$

Nonlinear dependency in data

Solution: basis functions (polynomials, gaussian, logistic sigmoid)

$$f_w(x) = w_0 + \sum_{j=1}^N w_j x^j \xrightarrow{\text{more generally}} f_w(x) = w_0 + \sum_{j=1}^N w_j \phi_j(x)$$

Defining $\phi_0 = 1$ $f_w(x) = w^T \phi(x)$ The function f
is still linear in w
(despite not being
linear in x)

$$f_w(x) = w_0 + \sum_{j=1}^m w_j \phi_j(x) = w^T \phi(x)$$

$$E_{LS}(w) = \frac{1}{2} \sum_i^n (w^T \phi(x_i) - y_i)^2 = \frac{1}{2} (\Phi w - y)^T (\Phi w - y)$$

$$\Phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_m(x_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_n) & \phi_1(x_n) & \dots & \phi_m(x_n) \end{pmatrix} \in \mathbb{R}^{n \times (m+1)}$$

$$\text{Optimal weights } w^* : w^* = (\Phi^T \Phi)^{-1} \Phi^T y = \Phi^+ y$$

To control **overfitting** \rightarrow we introduce regularization

$$E_{\text{ridge}}(w) = \underbrace{\frac{1}{2} \sum_{i=1}^n [w^T \phi(x_i) - y_i]^2}_{\text{large values}} + \underbrace{\frac{\lambda}{2} \|w\|_2^2}_{\text{low values}}$$

- $\|w\|_2^2 = w^T w = w_0^2 + w_1^2 + \dots + w_M^2$ - squared L₂ norm of w
- $\lambda \rightarrow$ regularization strength

The error of an estimator can be decomposed into 2 parts:

- Bias \rightarrow expected error due to model mismatch
- Variance \rightarrow variation due to randomness in training data

High bias \rightarrow the model is too rigid to fit the underlying data distribution
(regularization strength λ is too high)

High variance \rightarrow the model is too flexible, and therefore captures noise in the data
OVERFITTING (high capacity, memorizes data and λ too low)

WHAT IS IDEAL \rightarrow low bias and low variance

The weights w_i can be interpreted as the strength of the (linear) relationship between feature x_i and y

But correlation does not imply causation

Probabilistic Linear Regression

$$y_i = f_w(x_i) + \varepsilon_i \quad \varepsilon_i \sim N(0, \beta^{-1})$$

(noise)

$$y_i \sim N(f_w(x_i), \beta^{-1})$$

maximum likelihood

of a single sample $\rightarrow p(y_i | f_w(x_i), \beta) = N(y_i | f_w(x_i), \beta^{-1})$

entire dataset $D = \{X, y\} \rightarrow p(y | X, w, \beta) = \prod_{i=1}^N p(y_i | f_w(x_i), \beta)$

$$w_{ML}, \beta_{ML} = \underset{w, \beta}{\operatorname{argmax}} p(y | X, w, \beta) = \underset{w, \beta}{\operatorname{argmax}} \ln p(y | X, w, \beta) =$$

$$\underset{w, \beta}{\operatorname{argmin}} -\ln p(y | X, w, \beta)$$

maximum likelihood error function $E_{ML}(w, \beta) = -\ln p(y | X, w, \beta)$

$$E_{ML}(w, \beta) = -\ln \left[\prod_{i=1}^N p(y_i | f_w(x_i), \beta) \right]$$

$$= -\ln \left[\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (w^\top \phi(x_i) - y_i)^2 \right) \right]$$

$$= -\sum_{i=1}^N \ln \left[\sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (w^\top \phi(x_i) - y_i)^2 \right) \right]$$

$$= \frac{\beta}{2} \sum_{i=1}^N (w^\top \phi(x_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi$$

Optimizing log-likelihood

$$w_{ML} = \underset{w}{\operatorname{argmin}} E_{ML}(w, \beta) = \underset{w}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^N (w^\top \phi(x_i) - y_i)^2 = \underset{w}{\operatorname{argmin}} E_{LS}(w)$$

Maximizing the likelihood is equivalent to minimizing the least squares error function

$$\beta_{ML} = \underset{\beta}{\operatorname{argmin}} E_{ML}(w_{ML}, \beta)$$

$$= \underset{\beta}{\operatorname{argmin}} \left[\frac{\beta}{2} \sum_{i=1}^N (w_{ML}^\top \phi(x_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \right]$$

take derivative β and set it to 0

$$\frac{\partial}{\partial \beta} E_{ML}(w_{ML}, \beta) = \frac{1}{2} \sum_{i=1}^N (w_{ML}^\top \phi(x_i) - y_i)^2 - \frac{N}{2\beta} \stackrel{!}{=} 0$$

$$\text{Solving for } \beta \quad \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{i=1}^N (w_{ML}^\top \phi(x_i) - y_i)^2$$

Posterior distribution

$$p(w | X, y, \beta, \cdot) = \frac{\underbrace{p(y | X, w, \beta)}_{\text{likelihood}} \cdot \underbrace{p(w | \cdot)}_{\text{prior}}}{\underbrace{p(X, y | w, \beta, \cdot)}_{\text{normalizing constant}}} \propto p(y | X, w, \beta) \cdot p(w | \cdot)$$

Prior for w

$$p(w | \alpha) = N(w | 0, \alpha^{-1} I) = \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp \left(-\frac{\alpha}{2} w^\top w \right)$$

$\alpha \rightarrow$ precision of the distribution
 $M \rightarrow n =$ number of elements in the vector w

Maximum a posteriori (MAP)

$$w_{MAP} = \underset{w}{\operatorname{argmax}} p(w | X, y, \alpha, \beta) = \underset{w}{\operatorname{argmax}} \ln p(y | X, w, \beta) + \ln p(w | \alpha) - \underbrace{\ln p(X, y | w, \beta, \alpha)}_{\text{const}}$$

$$= \underset{w}{\operatorname{argmin}} -\ln p(y | X, w, \beta) - \ln p(w | \alpha)$$

$$E_{MAP}(w) = -\ln p(y | X, w, \beta) - \ln p(w | \alpha) + \text{const}$$

$$= \frac{\beta}{2} \sum_{i=1}^N (w^\top \phi(x_i) - y_i)^2 + \frac{\alpha}{2} \|w\|_2^2 + \text{const} \propto \text{Eridge}(w) + \text{const}$$

MAP estimation with Gaussian prior β equivalent to ridge regression

Full Bayesian approach

The full posterior distribution $p(w|\mathcal{D}) \propto p(y|X, w, \beta) \cdot p(w|\alpha)$

Since both likelihood and prior are Gaussian, the posterior is as well

$$p(w|\mathcal{D}) = N(w|\mu, \Sigma)$$

$$\mu = \beta \Sigma \Phi^T y \quad \| \quad \Sigma^{-1} = \alpha I + \beta \Phi^T \Phi$$

$$w_{MAP} = \mu \quad \text{for } N=0, \text{ posterior} = \text{prior}$$

$$\alpha \rightarrow 0 \Rightarrow w_{MAP} = w_{ML}$$

Predictive distribution

\hat{y}_{new} for new data x_{new}

- Maximum Likelihood: w_{ML} and β_{ML}

$$p(\hat{y}_{\text{new}} | x_{\text{new}}, w_{ML}, \beta_{ML}) = N(\hat{y}_{\text{new}} | w_{ML}^\top \phi(x_{\text{new}}), \beta_{ML}^{-1})$$

- Maximum a posteriori: w_{MAP}

$$p(\hat{y}_{\text{new}} | x_{\text{new}}, w_{MAP}, \beta) = N(\hat{y}_{\text{new}} | w_{MAP}^\top \phi(x_{\text{new}}), \beta^{-1})$$

posterior predictive distribution

$$p(\hat{y}_{\text{new}} | x_{\text{new}}, \mathcal{D}) = N(\hat{y}_{\text{new}} | \mu^\top \phi(x_{\text{new}}), \beta^{-1} + \phi(x_{\text{new}})^\top \Sigma \phi(x_{\text{new}}))$$