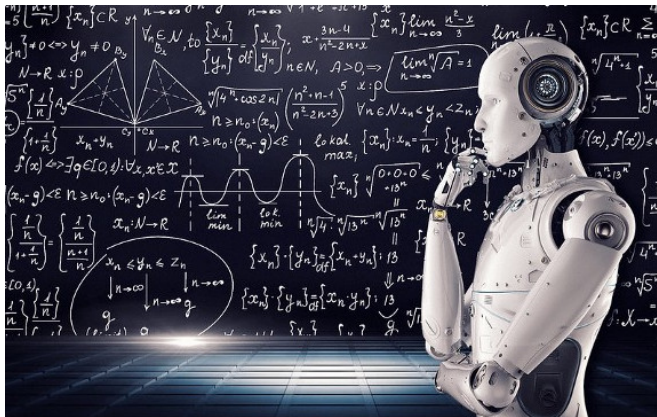


Common Machine Learning Algorithms



CONTENTS

1 Linear Models

LINEAR MODELS

LINEAR MODELS – FUNCTIONALITY

SUPERVISED

PARAMETRIC

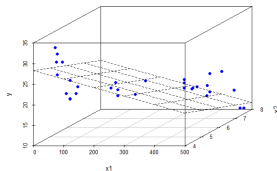
WHITE-BOX

General idea Represent target as function of linear predictor $\theta^T \mathbf{x}$

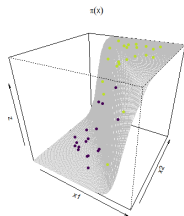
Hypothesis space

$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathbb{R} \mid f(\mathbf{x}) = \phi(\theta^T \mathbf{x})\}$, with suitable transformation $\phi(\cdot)$, e.g.,

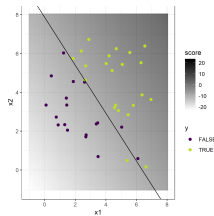
- Identity $\phi(\theta^T \mathbf{x}) = \theta^T \mathbf{x} \Rightarrow$ **linear regression**
- Logistic sigmoid function $\phi(\theta^T \mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})} =: \pi(\mathbf{x} \mid \theta) \Rightarrow$ **(binary) logistic regression**
 - Probability $\pi(\mathbf{x} \mid \theta) = \mathbb{P}(y = 1 \mid \mathbf{x})$ of belonging to one of two classes
 - Separating hyperplane via decision rule (e.g., $\hat{y} = 1 \Leftrightarrow \pi(\mathbf{x}) > 0.5$)



Linear regression hyperplane



Logistic function for bivariate input



Separating hyperplane for bivariate logistic regression

LINEAR MODELS – FUNCTIONALITY

Empirical risk

- **Linear regression**

- Typically, based on **quadratic** loss: $\mathcal{R}_{\text{emp}}(\theta) = \sum_{i=1}^n \left(y^{(i)} - f(\mathbf{x}^{(i)} | \theta) \right)^2$
⇒ corresponding to ordinary-least-squares (OLS) estimation
- Alternatives: e.g., **absolute** or **Huber** loss (both improving robustness)

- **Logistic regression:** based on **Bernoulli/log/cross-entropy** loss

$$\Rightarrow \mathcal{R}_{\text{emp}}(\theta) = \sum_{i=1}^n -y^{(i)} \log \left(\pi(\mathbf{x}^{(i)}) \right) - (1 - y^{(i)}) \log \left(1 - \pi(\mathbf{x}^{(i)}) \right)$$

Optimization

- For **OLS**: analytically with $\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
(with $\mathbf{X} \in \mathbb{R}^{n \times p}$: matrix of feature vectors)
- For **other loss functions**: numerical optimization

Hyperparameters None

LINEAR MODELS – PRO'S & CON'S

Advantages

- + **simple and fast** implementation
- + **analytical** solution
- + **cheap** computation
- + applicable for any **dataset size**, as long as number of observations \gg number of features
- + flexibility **beyond linearity** with polynomials, trigonometric transformations etc.
- + intuitive **interpretability** via feature effects
- + statistical hypothesis **tests** for effects available

Disadvantages

- **nonlinearity** of many real-world problems
- further restrictive **assumptions**: linearly independent features, homoskedastic residuals, normality of conditional response
- **sensitivity** w.r.t. outliers and noisy data (especially with L2 loss)
- risk of **overfitting** with more flexible feature transformations (e.g., high-degree polynomials)
- feature **interactions** must be handcrafted, so higher orders practically infeasible
- no handling of **missing** data

Simple method with good interpretability for linear problems, but with strong assumptions and limited complexity

LINEAR MODELS – PRACTICAL HINTS

Implementation

- **R:** `mlr3 learner LearnerRegrLM`, calling `stats::lm()` / `mlr3 learner LearnerClassifLogReg`, calling `stats::glm()`
- **Python:** `LinearRegression` from package `sklearn.linear_model`, package for advanced statistical parameters `statsmodels.api`

Regularization

- Quadratic penalty on model coefficients (a.k.a. **Ridge** regression)
 - Overall smaller, but still dense θ
 - Suitable with many influential features present, handling correlated features by shrinking their coefficients equally
- Absolute penalty on model coefficients (a.k.a. **Lasso** regression)
 - Actual variable selection
 - Suitable for sparse problems, ineffective with correlated features (randomly selecting one)
- Weighted combination of Ridge and Lasso: **elastic net**