

A TOUR OF MACHINE LEARNING CLASSIFIERS USING SCIKIT-LEARN

- NO single classifier works best across all possible scenarios
- in practice, it is recommended that you compare the performance of at least a handful of different algorithms & select the best model for the particular problem.

important factors:

- # of features
- amount of noise
- whether classes are linearly separable.

• 5 main steps for training:

1. feature selection & collection of labeled training data.
2. choosing a performance metric (MSE vs R²)
3. choosing an algorithm & training the model
4. Evaluate performance of the model
5. tune model performance.

- train/test split from sklearn pre-shuffles the dataset.
- **Stratification**: train/test split will return splits that have the same proportions of class labels as the input dataset

***also called z-score normalization**

- **Standard Scaler**: standard scaler uses the fit method to estimate the mean & standard deviation of each feature, where transform actually makes features have a mean of zero & a standard dev of 1.
- **one-vs-rest**: the model handles multi-classification tasks by creating n binary classifiers where n is the # of unique classes. that way for the given class, the the desired class is given positive outcomes, rest are given negative.

• **Error vs. Accuracy**

$$\text{Error} = \frac{\# \text{ of incorrect predictions}}{\# \text{ of predictions}}$$
$$\text{Accuracy} = (1 - \text{Error}) \times 100 (\%)$$

- **score method**: using the score method combines the predict/accuracy score call

- perceptron will never converge if the classes are not linearly separable

LOGISTIC REGRESSION

- LR works based on defining odds & the probability of a "positive event"

ODDS = $p/(1-p)$, where p = probability of a positive event

$$p := p(y=1 | x)$$

label \swarrow \nwarrow features
 natural log \nearrow

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

- under a logistic model we assume a linear relationship between our inputs & log-odds:

$$\text{logit}(p) = w_1x_1 + w_2x_2 + \dots + w_mx_m + b = \sum_{i=1}^m w_ix_i + b = \mathbf{w}^T \mathbf{x} + b$$

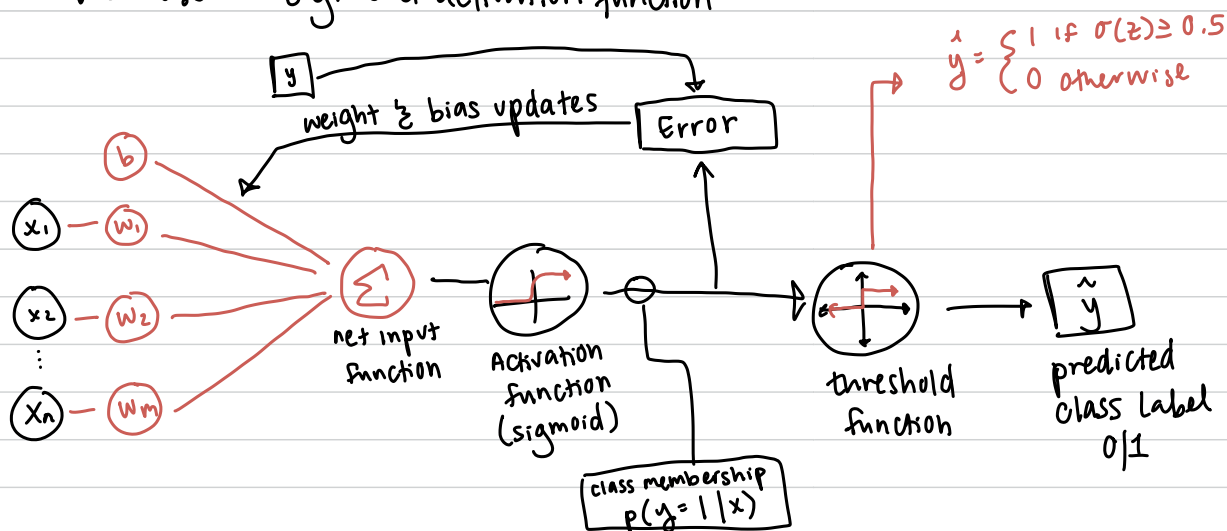
we are interested in solving for p .
 so we take in inverse nat. log to find $\sigma(z)$.

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \text{ where } z = \mathbf{w}^T \mathbf{x} + b$$

\nwarrow this is called the logistic sigmoid function.

if $z \rightarrow \infty$, $\sigma(z) = 1$ ($e^{-\infty} \rightarrow 0$). Similarly, if $z \rightarrow -\infty$, $\sigma(z) = 0$

- this works similar Adaline, where instead of an identity function, we use a sigmoid activation function



LEARNING MODEL WEIGHTS VIA LOGISTIC LOSS FUNCTION

ADALINE MSE loss:

$$L(w, b | \mathcal{X}) = \sum_i \frac{1}{2} (\sigma(z^i) - y^i)^2$$

$$p(Y=1 | x=x^i) = \sigma(z^i)$$

$$p(Y=0 | x=x^i) = 1 - \sigma(z^i)$$

- Assuming our examples are independent, we can define the likelihood \mathcal{L} that we want to maximize in logistic regression: Prob. mass func.

$$L(w, b | \mathcal{X}) = p(y | \mathcal{X}; w, b) = \prod_{i=1}^n p(y^i | x^i; w, b) = \prod_{i=1}^n (\sigma(z^i))^{y^i} \cdot (1 - \sigma(z^i))^{1-y^i}$$

Bernoulli principle

- It's much easier to maximize the nat. log of this function, which is called the log-likelihood:

$$\ell(w, b | \mathcal{X}) = \log \mathcal{L}(w, b | \mathcal{X}) = \sum_{i=1}^n [y^i \log(\sigma(z^i)) + (1-y^i) \log(1-\sigma(z^i))]$$

* we can maximize this using gradient ascent *

For gradient descent:

$$\mathcal{L}(w, b) = \sum_{i=1}^n [-y^i \log(\sigma(z^i)) - (1-y^i) \log(1-\sigma(z^i))]$$

$$\therefore \mathcal{L}(\sigma(z), y; w, b) = -y \log(\sigma(z)) - (1-y) \log(1-\sigma(z))$$

$$\mathcal{L}(\sigma(z), y; w, b) = \begin{cases} -\log(\sigma(z)) & \text{if } y=1 \\ -\log(1-\sigma(z)) & \text{if } y=0 \end{cases}$$

↑ maximize = Add
↓ minimize = Subtract

TACKLING OVERFITTING VIA REGULARIZATION

- overfitting means that the model is too well fit for the training data, & doesn't generalize well to unseen data.
- If the model is overfit, the model will have high variance, which is caused by having too many parameters (too complex for the data)
- If the model is underfit, it will have high bias, meaning the data is too complex for the model

* this is called the bias/variance tradeoff *

- to find a bias-variance tradeoff, we can tune the complexity of the model via Regularization.
- Regularization is very useful method for handling collinearity (large correlation between features), filtering out noise, & preventing overfitting.

eigenvector ↗

$$\text{L2 regularization: } \frac{\lambda}{2n} \|W\|^2 = \frac{\lambda}{2n} \sum_{j=1}^m w_j^2$$

HIGHLY SENSITIVE TO SCALING

loss changes to:

$$\mathcal{L}(w, b) = \frac{1}{n} \sum_{i=1}^n [-y^i \log(\sigma(z^i)) - (1 - y^i) \log(1 - \sigma(z^i))] + \frac{\lambda}{2n} \|W\|^2$$

- Regularization is useful to prevent overfitting, but if the strength is too great, it will also lead to underfitting.