# Machine Learning Tips

## Metrics

Given a set of data points $\{x^{(1)}, ..., x^{(m)}\}$, where each $x^{(i)}$ has $n$ features, associated to a set of outcomes $\{y^{(1)}, ..., y^{(m)}\}$, we want to assess a given classifier that learns how to predict $y$ from $x$.

## Classification

In a context of a binary classification, here are the main metrics that are important to track to assess the performance of the model.

❐ **Confusion matrix** – The confusion matrix is used to have a more complete picture when assessing the performance of a model. It is defined as follows:

**Predicted** class

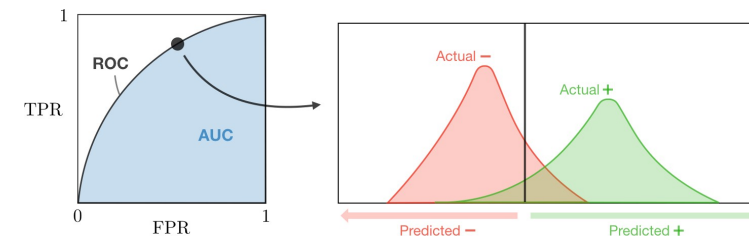|  | + | − |
|---|---|---|
| **+** | **TP** True Positives | **FN** False Negatives Type II error |
| **−** | **FP** False Positives Type I error | **TN** True Negatives |

**Actual** class

❐ **Main metrics** – The following metrics are commonly used to assess the performance of classification models:

| Metric | Formula | Interpretation |
|---|---|---|
| Accuracy | $\dfrac{TP+TN}{TP+TN+FP+FN}$ | Overall performance of model |
| Precision | $\dfrac{TP}{TP+FP}$ | How accurate the positive predictions are |
| Recall Sensitivity | $\dfrac{TP}{TP+FN}$ | Coverage of actual positive sample |
| Specificity | $\dfrac{TN}{TN+FP}$ | Coverage of actual negative sample |
| F1 score | $\dfrac{2TP}{2TP+FP+FN}$ | Hybrid metric useful for unbalanced classes |

❐ **ROC** – The receiver operating curve, also noted ROC, is the plot of TPR versus FPR by varying the threshold. These metrics are are summed up in the table below:

| Metric | Formula | Equivalent |
|---|---|---|
| True Positive Rate TPR | $\dfrac{TP}{TP+FN}$ | Recall, sensitivity |
| False Positive Rate FPR | $\dfrac{FP}{TN+FP}$ | 1-specificity |

❐ **AUC** – The area under the receiving operating curve, also noted AUC or AUROC, is the area below the ROC as shown in the following figure:



## Regression

❐ **Basic metrics** – Given a regression model $f$, the following metrics are commonly used to assess the performance of the model:

| Total sum of squares | Explained sum of squares | Residual sum of squares |
|---|---|---|
| $SS_{tot}= \displaystyle\sum_{i=1} (y_i - \bar{y})^2$ | $SS_{reg}= \displaystyle\sum_{i=1} (f(x_i) - y)^2$ | $SS_{res}= \displaystyle\sum_{i=1} (y_i - f(x_i))^2$ |

❐ **Coefficient of determination** – The coefficient of determination, often noted $R^2$ or $r^2$, provides a measure of how well the observed outcomes are replicated by the model and is defined as follows:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

❐ **Main metrics** – The following metrics are commonly used to assess the performance of regression models, by taking into account the number of variables $n$ that they take into consideration:

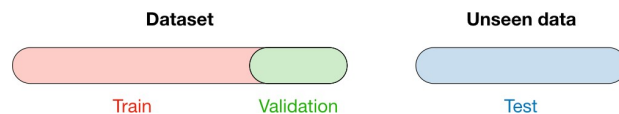| Mallow's Cp | AIC | BIC | Adjusted $R^2$ |
|---|---|---|---|
| $\dfrac{SS_{res} + 2(n+1)\hat{\sigma}^2}{}$ | $2\left((n+2) - \log(L)\right)$ | $\log(m)(n+2) - 2\log(L)$ | $1 - \dfrac{(1 - R^2)(m - 1)}{m - n - 1}$ |

where $L$ is the likelihood and $\sigma^2$ is an estimate of the variance associated with each response.

## Model selection

❒ **Vocabulary** – When selecting a model, we distinguish 3 different parts of the data that we have as follows:

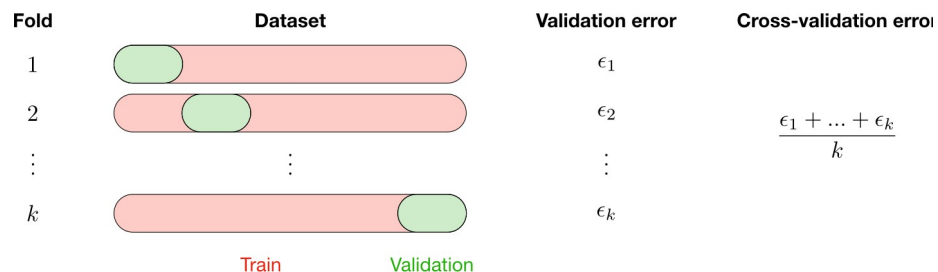| Training set | Validation set | Testing set |
|---|---|---|
| - Model is trained<br>- Usually 80% of the dataset | - Model is assessed<br>- Usually 20% of the dataset<br>- Also called hold-out or development set | - Model gives predictions<br>- Unseen data |

Once the model has been chosen, it is trained on the entire dataset and tested on the unseen test set. These are represented in the figure below:

**Dataset**          **Unseen data**

Train    Validation          Test

❒ **Cross-validation** – Cross-validation, also noted CV, is a method that is used to select a model that does not rely too much on the initial training set. The different types are summed up in the table below:

| $k$-fold | Leave-$p$-out |
|---|---|
| - Training on $k-1$ folds and assessment on the remaining one<br>- Generally $k = 5$ or $10$ | - Training on $n-p$ observations and assessment on the $p$ remaining ones<br>- Case $p = 1$ is called leave-one-out |

The most commonly used method is called $k$-fold cross-validation and splits the training data into $k$ folds to validate the model on one fold while training the model on the $k-1$ other folds, all of this $k$ times. The error is then averaged over the $k$ folds and is named cross-validation error.

| Fold | Dataset | Validation error | Cross-validation error |
|---|---|---|---|
| 1 | | $\epsilon_1$ | |
| 2 | | $\epsilon_2$ | $\dfrac{\epsilon_1 + \dots + \epsilon_k}{k}$ |
| ⋮ | ⋮ | ⋮ | |
| $k$ | | $\epsilon_k$ | |

Train          Validation

❒ **Regularization** – The regularization procedure aims at avoiding the model to overfit the data and thus deals with high variance issues. The following table sums up the different types of commonly used regularization techniques:

| LASSO | Ridge | Elastic Net |
|---|---|---|
| - Shrinks coefficients to 0<br>- Good for variable selection | Makes coefficients smaller | Tradeoff between variable selection and small coefficients |
| $\|\|\theta\|\|_1 \leqslant 1$ | $\|\|\theta\|\|_2 \leqslant 1$ | $(1-\alpha)\|\|\theta\|\|_1 + \alpha\|\|\theta\|\|_2^2 \leqslant 1$ |
| $\dots + \lambda\|\|\theta\|\|_1$<br>$\lambda \in \mathbb{R}$ | $\dots + \lambda\|\|\theta\|\|_2^2$<br>$\lambda \in \mathbb{R}$ | $\dots + \lambda\left[(1-\alpha)\|\|\theta\|\|_1 + \alpha\|\|\theta\|\|_2^2\right]$<br>$\lambda \in \mathbb{R}, \alpha \in [0,1]$ |

❒ **Model selection** – Train model on training set, then evaluate on the development set, then pick best performance model on the development set, and retrain all of that model on the whole training set.

## Diagnostics

❒ **Bias** – The bias of a model is the difference between the expected prediction and the correct model that we try to predict for given data points.

❒ **Variance** – The variance of a model is the variability of the model prediction for given data points.

❒ **Bias/variance tradeoff** – The simpler the model, the higher the bias, and the more complex the model, the higher the variance.

| | Underfitting | Just right | Overfitting |
|---|---|---|---|
| **Symptoms** | - High training error<br>- Training error close to test error<br>- High bias | - Training error slightly lower than test error | - Low training error<br>- Training error much lower than test error<br>- High variance |
| **Regression** | | | |

| | | | |
|---|---|---|---|
| **Classification** |  |  |  |
| **Deep learning** |  |  |  |
| **Remedies** | - Complexify model<br>- Add more features<br>- Train longer | | - Regularize<br>- Get more data |

❑ **Error analysis**– Error analysis is analyzing the root cause of the difference in performance between the current and the perfect models.

❑ **Ablative analysis** – Ablative analysis is analyzing the root cause of the difference in perfor- mance between the current and the baseline models.