

Supervised Regression

$$H = \{f(x) = \theta^T x \mid \theta \in \mathbb{R}^{p+1}\}$$

L = L2-Loss

$$R_{\text{emp}}(\theta) = \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^p} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

$$= \arg \min_{\theta \in \mathbb{R}^p} \|y - X\theta\|_2^2 = (X^T X)^{-1} X^T y$$

Overfitting and Underfitting

Underfitting occurs when a model can't reflect

the true shape of underlying function (given the data)

↳ High Test Error and high Training Error.

Overfitting occurs when a model reflects noise

or artifacts from D_{train} which do not generalize.

↳ Small Train Error and high Test Error.

Overfitting is influenced by:

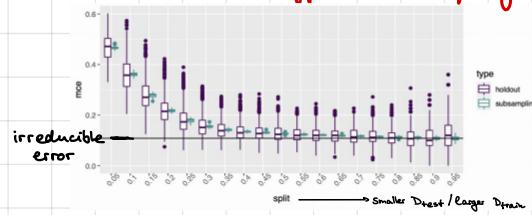
Complexity $\dim(\theta)$, n_{train} , $\dim(x)$, aleatoric uncertainty

Overfitting can be avoided by regularization:

less complex model, quality of training data,

Feature Selection/Engineering, Occam's razor

Bias-Variance-Tradeoff in Subsampling



$$\text{Bias}(\text{Subsampling} \mid \text{split-rate}) = \text{Bias}(\text{Hold-out} \mid \text{split-rate})$$

$$\text{Variance}(\text{Subsampling} \mid \text{split-rate}) < \text{Variance}(\text{Hold-out} \mid \text{split-rate})$$

⇒ "Optimal" split-rate is higher in Subsampling compared to one Hold-out-split.

Bias-Variance Analysis in Resampling

If there exists a dedicated Test set, then

$$\widehat{GE}(\hat{f}, L) = \frac{1}{m} \sum_{(x_i, y_i) \in D_{\text{test}}} L(y_i, \hat{f}(x_i))$$

$$\mathbb{E}[\widehat{GE}(\hat{f}, L)] = GE(\hat{f}, L) \text{ and } V[\widehat{GE}(\hat{f}, L)] = \frac{1}{m} V[L(y_i, \hat{f}(x_i))]$$

We can apply the CLT to approx. the distr. of $\widehat{GE}(\hat{f}, L)$.

$$\mathbb{E}[\widehat{GE}(I, J, p)] \approx GE(I, n_{\text{train}}, p) \leq GE(I, n, p)$$

Supervised Classification

Scoring classifiers: scoring/discriminant functions $f_1, \dots, f_g : \mathcal{X} \rightarrow \mathbb{R}$.

$$\text{classifier } h(x) := \arg \max_{k=1, \dots, g} f_k(x)$$

Prob. Classifier: probability functions $\pi_1, \dots, \pi_g : \mathcal{X} \rightarrow [0, 1]$ with $\sum_{k=1}^g \pi_k = 1$.

$$\text{classifier } h(x) := \arg \max_{k=1, \dots, g} \pi_k(x)$$

Linear Classifier: f_1, \dots, f_g are called linear specifiers if

$$\exists g \text{ rank-preserving, monotone: } g(f_k(x)) = w_k^T x + b_k$$

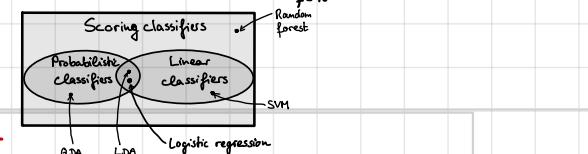
Approaches to construct classifiers

Generative Approach: "Which y tends to have x like that?"

$$\pi_k(x) \propto P(x \mid y=k) \cdot \pi_k \quad \pi_k \text{ being a prior}$$

Discriminant Approach "What is the best prediction for y given x ?"

$$\hat{f} = \arg \min_{f \in H} R_{\text{emp}}(f)$$



Train Error vs. Test Error

Goodness-of-fit measures - like R^2 , χ^2 , AIC, BIC, deviance - is based on training error but are based on distributional assumptions and limited data - therefore hard to use for high-dimensional or more complex data.

Decrease of n_{train} → Increase of Test Error (in general) b.c. model generalizes better with more training data and worse with less training data.
 Increase of complexity → Decrease of Training Error (in general) b.c. it becomes easier to learn all patterns on small training data sets or with more flexibility.
 Decrease of n_{test} → Increase of Variance of test error.

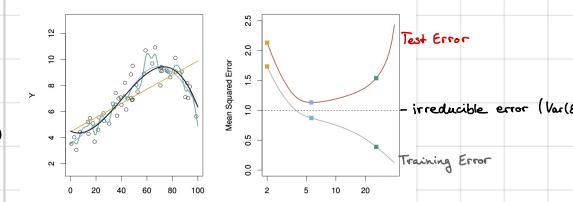
Bias-Variance - Tradeoff

Influence of n_{train} :

Because in practice $n = |\mathcal{D}|$ is fixed increase/decrease of n_{train} results in a decrease/increase of n_{test} .

Increasing n_{train} leads to decrease in Test Error but higher variance of Test Error.

Influence of complexity:



Examples

Logistic Regression (Discriminant):

$$H = \{ \pi : \mathcal{X} \rightarrow [0, 1] \mid \pi(x) = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)} \}$$

Discriminant Analysis (Generative)

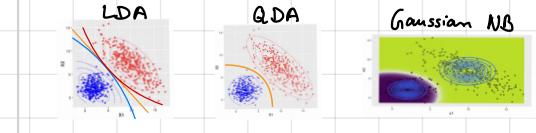
- LDA: assumes $x \mid y=k \sim N(\mu_k, \Sigma)$ (Linear)

- QDA: assumes $x \mid y=k \sim N(\mu_k, \Sigma_k)$

Naive Bayes (Generative) assume $p(x_j \mid y=k) = \prod_{j=1}^p p(x_j \mid y=k)$

Gaussian NB: assume $x_j \mid y=k \sim N(\mu_{jk}, \sigma_{jk}^2)$. Σ_{jk} diagonal matrix.

Multinomial NB: assume $p(x_j \mid y=k) \propto \prod_{m=1}^M p_{jm}$. P_{jm} = rel. freq. of m in feature j restricted to class k



ROC Analysis

		True condition		Prevalence = $\frac{\text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\text{True positive} + \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Total population	Condition positive	Condition negative		
	Predicted condition positive	True positive, Power	False positive, Type I error	True negative	False negative, Type II error
Predicted condition	Predicted condition negative	False negative, Type II error	True negative	False predictive value (PPV), Precision = $\frac{\text{True positive}}{\text{True positive} + \text{False positive}}$	False omission rate (FOR) = $\frac{\text{False negative}}{\text{True negative} + \text{False negative}}$
		True positive rate (TPR), Recall, Sensitivity, probability of detection = $\frac{\text{True positive}}{\text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\text{False positive}}{\text{Condition positive}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Negative predictive value (NPV) = $\frac{\text{True negative}}{\text{True negative} + \text{False negative}}$
		False negative rate (FNR), Miss rate = $\frac{\text{False negative}}{\text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\text{True negative}}{\text{Condition positive}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR}+}{\text{LR}-} = \frac{\text{ENR}}{\text{TNR}}$
				Negative likelihood ratio (LR-) = $\frac{\text{ENR}}{\text{TNR}}$	F1 score = $\frac{1}{\frac{\text{Recall}}{2} + \frac{\text{Precision}}{2}} = \frac{2 \cdot \text{TP}}{2 \cdot \text{TP} + \text{FN} + \text{FP}}$

→ Not a good performance measure when label dist. is unbalanced.

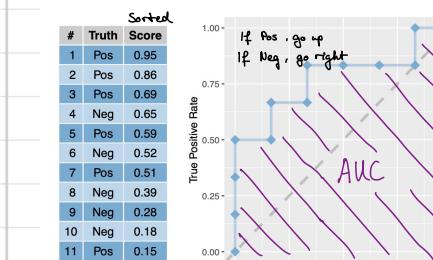
Accuracy paradox

Usually preferable in situations with imbalanced data.

→ Balances TPR and PPV but tends to the smaller value.

Different metrics emphasize different aspects of performance. Choice requires Domain Knowledge.

DRAWING ROC CURVES

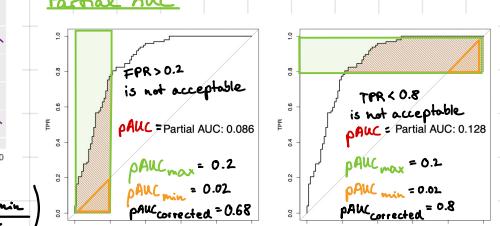


$$\text{PAUC}_{\text{corrected}} = \frac{1}{2} \left(1 + \frac{\text{PAUC}_{\text{max}} - \text{PAUC}_{\text{min}}}{\text{PAUC}_{\text{max}} - \text{PAUC}_{\text{min}}} \right)$$

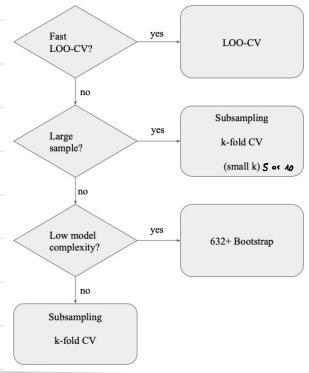
The scores are given by the classifier for each observation. Step-size is determined by # Pos and # Neg

Total

Partial AUC



Guidelines



Measures for Regression

Pointwise outer-Loss

Mean squared error (MSE): $\rho_{\text{MSE}}(\hat{y}, F) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$

Mean absolute error (MAE): $\rho_{\text{MAE}}(\hat{y}, F) = \frac{1}{m} \sum_{i=1}^m |y^{(i)} - \hat{y}^{(i)}|$

Mean abs. perc. error (MAPE): $\rho_{\text{MAPE}}(\hat{y}, F) = \frac{1}{m} \sum_{i=1}^m \frac{|y^{(i)} - \hat{y}^{(i)}|}{y^{(i)}}$

Set-based outer-Loss

$$R^2: \rho_R(\hat{y}, F) = 1 - \frac{\rho_{\text{MSE}}(\hat{y}, LM(\hat{y}))}{\rho_{\text{MSE}}(\hat{y}, \bar{y})}$$

! Higher $R^2 \not\Rightarrow$ Better fit.

But R^2 is invariant w.r.t linear scaling of y . MSE is not.

Generalized R^2 :

$$1 - \frac{\text{Loss}_{\text{ComplexModel}}}{\text{Loss}_{\text{SimpleModel}}}$$

E.g. model vs. constant,
linear vs. non-linear,
tree vs. forest, ...

! Usually $R^2 \in [0, 1]$, but this is only true if we evaluate on Training Data. On Test Data $R^2 < 0$ is possible \rightarrow Overfitting

Proximities

- Imputing missing data:
 - Replace missing values for a given variable using the median of the non-missing values
 - Get proximities
 - Replace missing values in observation $x^{(i)}$ by a weighted average of non-missing values, with weights proportional to the proximity between observation $x^{(i)}$ and the observations with the non-missing values
- Steps 2 and 3 are then iterated a few times.

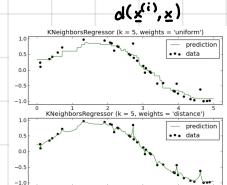
- Locating outliers:
 - An outlier is an observation whose proximities to all other observations are small
 - Measure of outlyingness can be computed for each observation in the training sample
 - If the measure is unusually large the observation should be carefully inspected
- Identifying mislabeled data:
 - Instances in the training data set are sometimes labeled ambiguously or incorrectly, especially in manually created data sets.
 - Proximities can help in finding them; they often show up as outliers in terms of their proximity values.
 - If the measure is unusually large the observation should be carefully inspected
- Visualizing the forest:
 - The values $1 - \text{prox}(x^{(i)}, x^{(j)})$ can be thought of as distances in a high-dimensional space
 - They can be projected onto a low-dimensional space using metric multidimensional scaling (MDS) eigenvectors of a modified version of the proximity matrix to get scaling coordinates

k-NN

Prediction Regression

$$\hat{f}(x) = \frac{1}{k} \sum_{i \in \mathcal{I}_k(x)} y^{(i)}$$

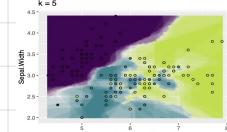
$$\hat{f}(x) = \frac{1}{\sum_{i \in \mathcal{I}_k(x)} w^{(i)}} \sum_{i \in \mathcal{I}_k(x)} w^{(i)} y^{(i)}$$



Prediction Classification

$$\hat{f}_k(x) = \frac{1}{k} \sum_{i \in \mathcal{I}_k(x)} \mathbb{1}\{y^{(i)} = l\}$$

$$\hat{h}(x) = \arg \max_{l \in \{1, \dots, g\}} \hat{f}_k(x)$$



Standardization and Weights

- Features are usually standardized or normalized.

because most distances place higher importance on features with higher ranges.

- Features can be given a higher importance by using a weighted distance measure.

$$d_{\text{Euclidean}}^{\text{Weighted}}(x, \hat{x}) = \sqrt{\sum_{j=1}^n w_j (x_j - \hat{x}_j)^2}$$

Measures for Classification

Pointwise outer-Loss using class-labels

Misclassification error rate (MCE):

$$\rho_{\text{MCE}} = \frac{1}{m} \sum_{i=1}^m [y^{(i)} \neq \hat{y}^{(i)}] = \frac{\sum \text{False pos. + False neg.}}{\text{Total population}}$$

Accuracy (ACC):

$$\rho_{\text{ACC}} = \frac{1}{m} \sum_{i=1}^m [y^{(i)} = \hat{y}^{(i)}] = \frac{\sum \text{True pos. + True neg.}}{\text{Total population}}$$

! No information about how good/skewed prob's are.
Errors on all classes weighted equally, which is often inappropriate

Cost matrix:

$$\begin{aligned} \text{Costs} &= \frac{1}{n} \sum_{i=1}^n C[y^{(i)}, \hat{y}^{(i)}] \\ &= \frac{1}{n} \langle \text{Cost Matrix}, \text{Confusion Matrix} \rangle_F \end{aligned}$$

Pointwise outer-Loss using probabilities

Brier Score

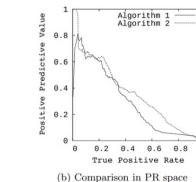
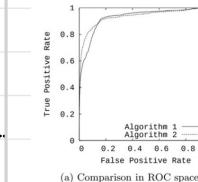
Multi-Class Brier Score

Log Loss

See Cheat-Sheet
! Take sum/mean

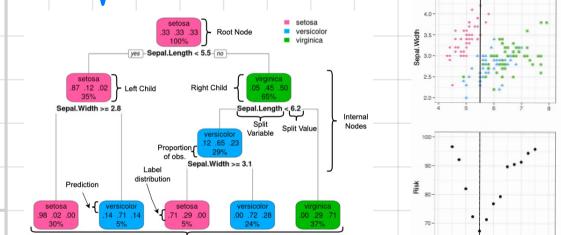
Precision - Recall Curves

Create a ROC-like plot, but with TPR and PPV instead of TPR and FPR. might be better for highly imbalanced data ($n_- \gg n_+$)

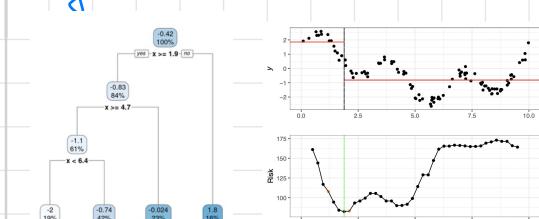


Curve dominates fully in ROC \Leftrightarrow
Curve dominates fully in PR.

Classification Trees



Regression Trees



Dealing with Categorical Features

If x_j is categorical with m levels, then there are $2^{m-1}-1$ possible partitions.

- binary classification:

◦ Calculate the proportion of 1-outcomes for each category of the feature in N .

◦ Sort the categories according to these proportions.

◦ The feature can then be treated as if it was ordinal, so we only have to investigate at most $m-1$ splits.

- regression with L2-loss:
 - Calculate the mean of the outcome in each category
 - Sort the categories by increasing mean of the outcome
 - The feature can then be treated as if it was ordinal, so we only have to investigate at most $m-1$ splits.

$$\widehat{e}_{\text{OOB}} = \frac{1}{n} \sum_{i=1}^n L(y^{(i)}, \hat{y}^{(i)}_{\text{OOB}})$$

$$\rho(e_{\text{OOB}}) = (1 - \frac{1}{n})^n \rightarrow \frac{1}{e} \approx 0.57$$

Neural Network
Identify an sigmoid function (Binary reg. or logistic reg.)

$$h = \mathbb{P} : \mathbb{R}^p \rightarrow \mathbb{R} \quad h(x) = \sigma\left(\sum_{i=1}^p w_i x_i + b\right), \text{ where } \sigma$$

As Loss we use $L(y, \hat{y})$ or $\text{Cross-Entropy Loss (log-reg.)}$