



Assessing Predictive Capability and Contribution for Binary Classification Models

Mindy Hotchkiss



Abstract



Classification models for binary outcomes are in widespread use across a variety of industries. Results are commonly summarized in a misclassification table, also known as an error or confusion matrix, which indicates correct vs incorrect predictions for different circumstances. Models are developed to minimize both false positive and false negative errors, but the optimization process to train/obtain the model fit necessarily results in cost-benefit trades. However, how to obtain an objective assessment of the performance of a given model in terms of predictive capability or benefit is less well understood, due to both the rich plethora of options described in literature as well as the largely overlooked influence of noise factors, specifically class imbalance. Many popular measures are susceptible to effects due to underlying differences in how the data are allocated by condition, which cannot be easily corrected.

This talk considers the wide landscape of possibilities from a statistical robustness perspective. Results are shown from sensitivity analyses for a variety of different conditions for several popular metrics and issues are highlighted, highlighting potential concerns with respect to machine learning or ML-enabled systems. Recommendations are provided to correct for imbalance effects, as well as how to conduct a simple statistical comparison that will detangle the beneficial effects of the model itself from those of imbalance. Results are generalizable across model type.

Overview



- Introduction & Terminology
- Motivating Example
- Issues & Concerns
- Implications for Machine Learning
- Candidate Metrics
- Sensitivity Study Results
- Alternative Approaches
- Summary & Recommendations

What are Classification Models?



- A mathematical function f used to predict outcome Y

$$f(X) \Rightarrow \hat{Y}$$

from inputs X

where Y is *NOT* continuous and takes on only a limited discrete set of values

- | |
|---|
| <ul style="list-style-type: none"> • Binary: 2 levels, denoting any dichotomy • Ordinal: 3+ classifications with natural sequencing • Nominal / Categorical: unordered categories or classifications |
|---|

- Standard regression models not applicable
 - Need a link function to convert to probability scale
- Results are summarized in an error matrix
 - Also designated contingency table, truth table, confusion matrix, misclassification table...

Raw Data Table

X1	...	Xk	Actual Y	Predicted Y
0.19	...	0.59	0	0
0.03	...	0.27	0	1
0.62	...	0.87	1	0
0.75	...	0.58	1	1
...



Binary Classification

		Model Prediction	
		No	Yes
Actual Condition	No	# True Negatives	# False Positives
	Yes	# False Negatives	# True Positives

Models must balance error types:

- **False Positive** – model is over-sensitive, detects too often
- **False Negative** – model is under-sensitive, fails to detect

How are Classification Models Assessed?



From this table, many measures can be calculated:

Binary Classification

		Model Prediction		Row Totals
		No	Yes	
Actual Condition	No	TN (# True Negatives)	FP (# False Positives)	AN (Total # Actual Negatives)
	Yes	FN (# False Negatives)	TP (# True Positives)	AP (Total # Actual Positives)
Col Totals		PN (Total # Predicted Negatives)	PP (Total # Predicted Positives)	N (Total)

Prevalence

$$= \frac{AP}{N}$$

* Class Imbalance *

Negative Predictive Value (NPV)

NPV

$$= \frac{TN}{PN} = \frac{TN}{FN + TN}$$

False Omission Rate (FOR = 1- NPV)

$$= \frac{FN}{PN} = \frac{FN}{FN + TN}$$

Positive Predictive Value (PPV)

$$= \frac{TP}{PP} = \frac{TP}{TP + FP}$$

Precision

False Discovery Rate (FDR = 1- PPV)

$$= \frac{FP}{PP} = \frac{FP}{TP + FP}$$

False Positive Rate (FPR = 1- TNR)

$$= \frac{FP}{AN} = \frac{FP}{FP + TN}$$

True Negative Rate (TNR)

$$= \frac{TN}{AN} = \frac{TN}{FP + TN}$$

True Positive Rate (TPR)

$$= \frac{TP}{AP} = \frac{TP}{TP + FN}$$

False Negative Rate (FNR = 1- TPR)

$$= \frac{FN}{AP} = \frac{FN}{TP + FN}$$

Fall-Out
Probability of
False Alarm

Specificity

Selectivity

Recall

Sensitivity

Hit Rate
Power

Probability of
Detection (POD)

Miss Rate

These are individual intermediate calculations, incomplete measures for model characterization or performance assessment

Introduction & Terminology

How are Classification Models Assessed?



(an incomplete list)

		Model Prediction	
		No	Yes
Actual Condition	No	TN (# True Negatives)	FP (# False Positives)
	Yes	FN (# False Negatives)	TP (# True Positives)

AN
(Total # Actual Negatives)

AP
(Total # Actual Positives)

N
(Total)

False Positive Rate (FPR = 1- TNR)

$$= \frac{FP}{AN} = \frac{FP}{FP+TN}$$

True Negative Rate (TNR)

$$= \frac{TN}{AN} = \frac{TN}{FP+TN}$$

True Positive Rate (TPR)

$$= \frac{TP}{AP} = \frac{TP}{TP+FN}$$

False Negative Rate (FNR = 1- TPR)

$$= \frac{FN}{AP} = \frac{FN}{TP+FN}$$

Fall-Out
Probability of False Alarm

Specificity
Selectivity

Recall
Sensitivity

Hit Rate
Power
Probability of Detection (POD)

Miss Rate

Functions of Row Proportions: Recall (Sensitivity) & Specificity

- Balanced Accuracy (BA) Arithmetic mean of Recall (Sensitivity) & Specificity
- Bookmaker Informedness (BMI), Youden's J

$$= \text{Recall} + \text{Specificity} - 1$$

$$= 2 * \text{BA} - 1 \text{ (Scaled BA)}$$

Functions of All Row & Column Proportions: Recall(Sensitivity), Specificity, NPV, & Precision

- Diagnostic Odds Ratio (DOR, OR), theta, cross-product ratio, Log DOR
- Yules Q, Y are scaled versions of theta
- Matthews Correlation Coefficient (MCC), Pearson's phi coefficient, also = Yules Y

Functions of Mixed Row & Column Proportions: Recall (Sensitivity) & Precision

- Area Under the Precision vs Recall Curve (AUC)
- Gini Coefficient = 2*AUC -1 (Scaled AUC)
- F-Score (F-Measure)

$$F_1: \text{Harmonic Mean of Precision \& Recall}$$

$$F_\beta: F_1 \text{ weighted to accommodate detection error trades}$$
- Fowlkes-Mallow Index: Geometric Mean of Precision & Recall

***** Omit True Negatives from all calculations*****

Functions of Neither Row nor Column Proportions

- Overall Accuracy

$$OA = \frac{TN+TP}{N}$$

Not all metrics are useful measures of model performance

Functions of Column Proportions: NPV & Precision

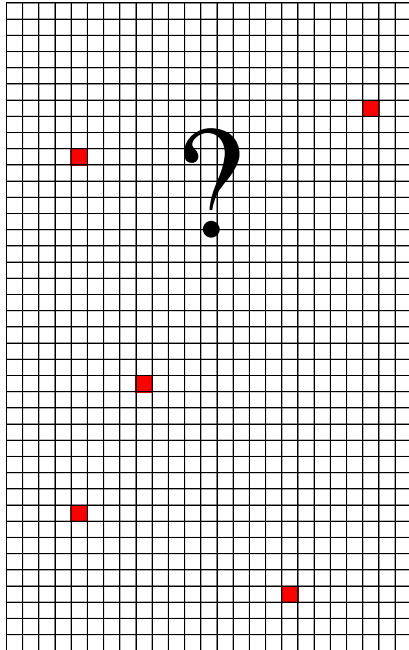
- Markedness (Mk) Δp

$$= \text{PPV(Precision)} + \text{NPV} - 1$$
 Calculated analogously to BMI & YJ except relative to columns instead of rows. Equivalently, could be scaled.

Rare Event Prediction



Example: 1000 predictions



- Objective is to successfully predict the outcome for each instance
- A high success rate can easily be achieved
... *without doing any actual modeling*
- Predict $Y = \text{No}$ for all instances (“null” or “naïve” model)

		Model Prediction		
		No	Yes	Total
Actual Condition	No	995	0	995
	Yes	5	0	5
	Total	1000	0	1000

$$\rightarrow \text{Overall Accuracy} = \frac{995+0}{1000} = \mathbf{99.5\%}$$

True incidence rate: 5/1000

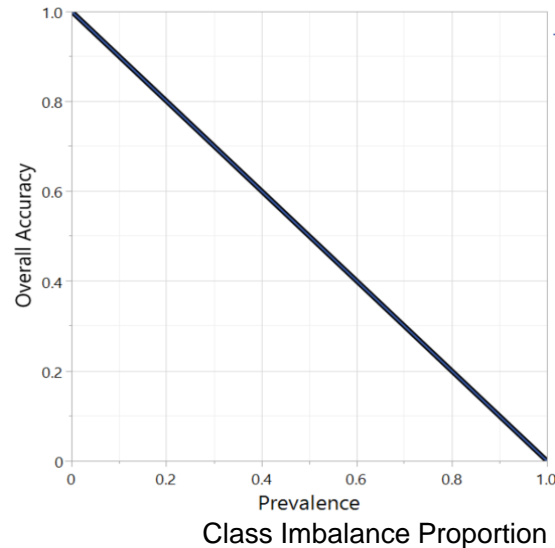
But this assessment is clearly meaningless – this “model” provides no information

Rare Event Prediction Implications



- If there is no value-added contribution from the model,

Overall Accuracy = 1 – Prevalence



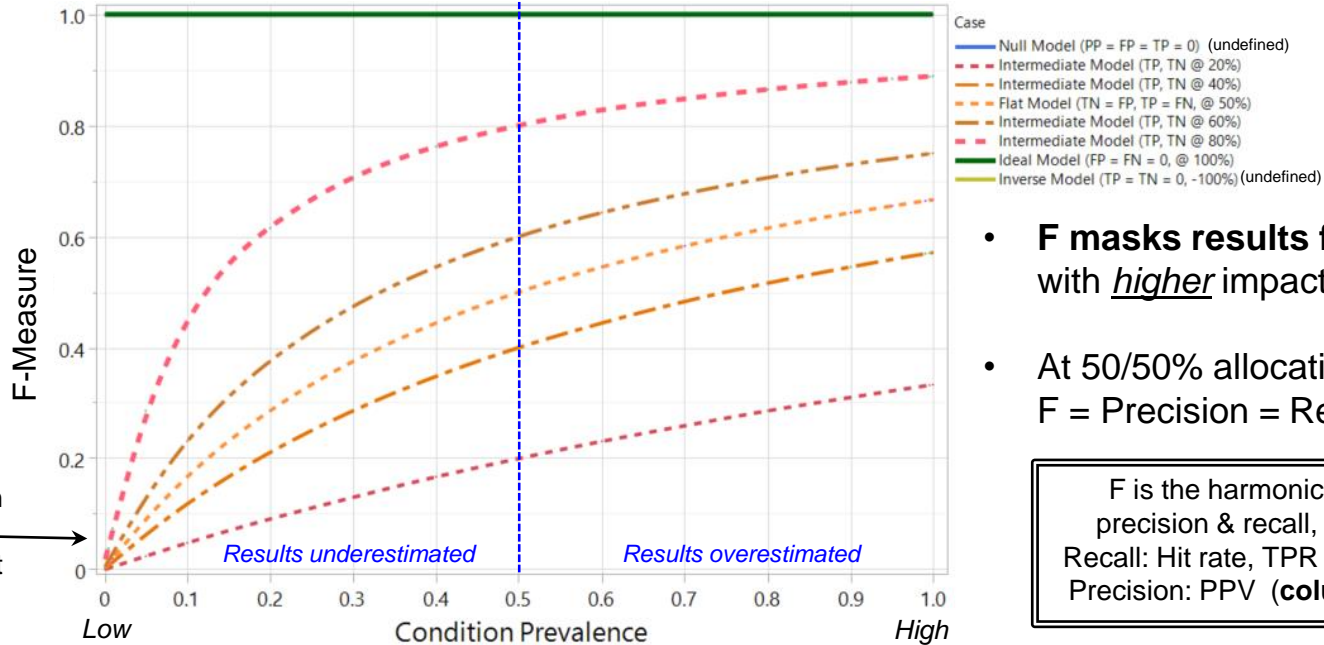
- If model effects are present, OA provides no way to disentangle its contribution from imbalance effects
- A good model will have high OA, but high OA itself is NOT an indicator of a good model
 - .. *necessary* but not *sufficient*
- OA is a *trailing* indicator, not a *leading* indicator

Overall Accuracy provides no insight into actual model performance



Class Imbalance Affects Other Metrics

- Sensitivity study results: F-measure also shows variable dependencies



- **F masks results for rarer events, with higher impact on better results**
- At 50/50% allocation, F = Precision = Recall

F is the harmonic mean of precision & recall, omits TN.
 Recall: Hit rate, TPR (**row** metric)
 Precision: PPV (**column** metric)

Undesirable properties generalizable to similarly derived functions

Robustness



- Statistical measures are tools with their own properties, that need to meet certain suitability criteria
- Tools that are more “robust” will be more broadly useful, less affected by violations in assumptions and usage conditions, but still providing a correct assessment
- Model assessment metrics should provide the SAME answer ***regardless of changes in “noise” factors that are unimportant***
- Class imbalance is unimportant, not integral to any analysis, and it should NOT drive conclusions

Implications for Machine Learning



- Class model assessment metrics in widespread use today have flaws that implicate results in majority of cases
- ML applications involving classification models should also have used a model assessment metric in development and during implementation if embedded into a larger more integrated system
- Down-stream decision-making based on these class model assessments assume that the metrics provide a fair and independent assessment of model performance when they do not
- Assessments of actual value contributed by classification models developed are highly likely to have been misrepresented
- Potential downstream effects include poor repeatability, increased inefficiencies, and loss of system credibility, as well as lost opportunities

Incorporating more robust class model metrics and assessments will provide greater chance of system success

Directionality affects Metric Robustness



Cross-Directional

		Model Prediction		
		No	Yes	Total
Actual Condition	No	True Negative (TN)	False Positive (FP)	Actual Negative (AN)
	Yes	False Negative (FN)	True Positive (TP)	Actual Positive (AP)
	Total	Predicted Negative (PN)	Predicted Positive (PP)	Total (N)

Metrics calculated relative to one column and one row, with some element omitted.

Unidirectional

		Model Prediction		
		No	Yes	Total
Actual Condition	No	True Negative (TN)	False Positive (FP)	Actual Negative (AN)
	Yes	False Negative (FN)	True Positive (TP)	Actual Positive (AP)
	Total	Predicted Negative (PN)	Predicted Positive (PP)	Total (N)

Metrics calculated relative to Columns

Unidirectional

		Model Prediction		
		No	Yes	Total
Actual Condition	No	True Negative (TN)	False Positive (FP)	Actual Negative (AN)
	Yes	False Negative (FN)	True Positive (TP)	Actual Positive (AP)
	Total	Predicted Negative (PN)	Predicted Positive (PP)	Total (N)

Metrics calculated relative to Rows

Bidirectional

		Model Prediction		
		No	Yes	Total
Actual Condition	No	True Negative (TN)	False Positive (FP)	Actual Negative (AN)
	Yes	False Negative (FN)	True Positive (TP)	Actual Positive (AP)
	Total	Predicted Negative (PN)	Predicted Positive (PP)	Total (N)

Metric calculated interchangeably relative to both rows AND columns

Balance is important: preferred metrics are in bidirectional family

Fully Balanced Functions



• Correlation: Phi, Matthews Correlation Coefficient (MCC)

- Range [-1,1], can be calculated on raw 0-1 data using standard Pearson correlation function
- Binary tabular formulation (one version):

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}$$

• (Diagnostic) Odds Ratio (DOR, OR), theta (θ), cross-product ratio

- Range [0, inf), 1 if no effect, undefined when FP, FN = 0
- Can be written equivalently in terms of *either* row or columns

$$OR = \frac{LR+}{LR-} = \frac{\frac{TPR}{FPR}}{\frac{FNR}{TNR}} = \frac{\frac{sensitivity}{1-specificity}}{\frac{1-sensitivity}{specificity}} = \frac{\frac{TP/AP}{FP/AP}}{\frac{TN/AN}{FN/AN}} = \frac{TP/FP}{FN/TN} = \frac{TP \cdot TN}{FP \cdot FN} = \frac{TP/FN}{FP/TN} = \frac{\frac{TP/PP}{FN/PN}}{\frac{FP/PP}{TN/PN}} = \frac{\frac{PPV}{FDR}}{\frac{NPV}{1-PPV}} = \frac{PPV}{NPV} = \theta$$

$$= \frac{sensitivity \cdot specificity}{(1 - sensitivity) \cdot (1 - specificity)} = \frac{PPV \cdot NPV}{(1 - PPV) \cdot (1 - NPV)}$$

– Variants Yule's Q, Y (coefficient of colligation)

- Functions of OR, converted to [-1,1] scale

$$Q = \frac{(\theta - 1)}{(\theta + 1)}; Y = \frac{(\sqrt{\theta} - 1)}{(\sqrt{\theta} + 1)}$$

*undefined when $\theta = -1$, NA since $\theta \in [0, \infty)$

Binary Case Layout

		Model Prediction		
		No	Yes	Total
Actual Condition	No	True Negative (TN)	False Positive (FP)	Actual Negative (AN)
	Yes	False Negative (FN)	True Positive (TP)	Actual Positive (AP)
Total		Predicted Negative (PN)	Predicted Positive (PP)	Total (N)

Metrics calculated using both rows AND columns

Functions of All Row & Column Proportions: Recall(Sensitivity), Specificity, NPV, & Precision

- Diagnostic Odds Ratio (DOR, OR), theta, cross-product ratio, Log DOR
- Yules Q, Y are scaled versions of theta
- Matthews Correlation Coefficient (MCC), Pearson's phi coefficient, also = Yules Y

Yule's Y conversion
= Delta Rate
= TNR-FPR

Sensitivity Study Comparisons & Results

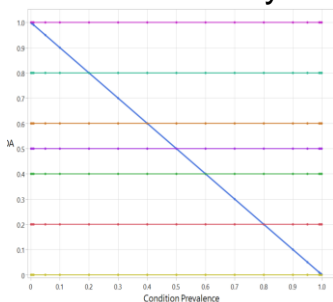


Case

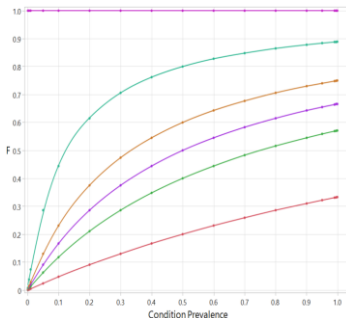
- Null Model (PP = FP = TP = 0)
- Intermediate Model (TP, TN @ 20%)
- Intermediate Model (TP, TN @ 40%)
- Flat Model (TN = FP, TP = FN, @ 50%)
- Intermediate Model (TP, TN @ 60%)
- Intermediate Model (TP, TN @ 80%)
- Ideal Model (FP = FN = 0, @ 100%)
- Inverse Model (TP = TN = 0, -100%)

**All cases defined to have specific OA*

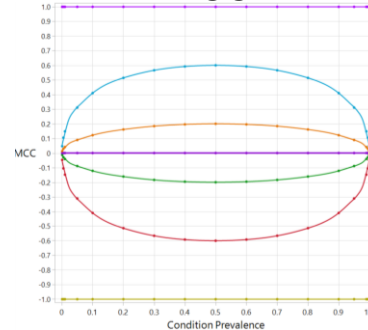
Overall Accuracy



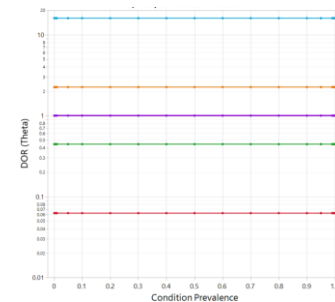
F-measure



MCC



DOR (logscale)

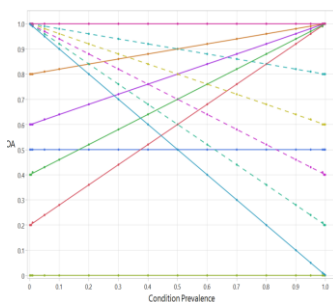


**Flat lines reflect good behavior*

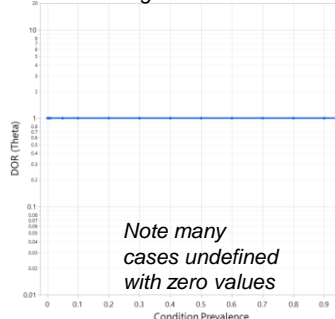
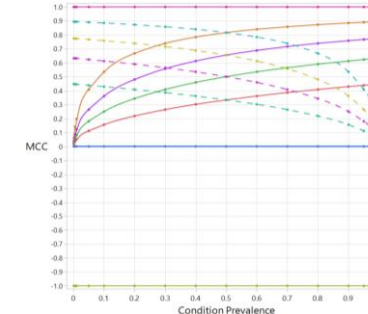
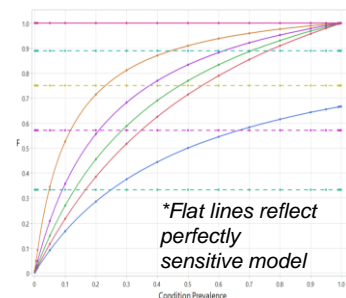
Case

- Flat Model (TN = FP, TP = FN, @ 50%)
- Highly Sensitive Model (FN = 0, TN @ 20%)
- Highly Sensitive Model (FN = 0, TN @ 40%)
- Highly Sensitive Model (FN = 0, TN @ 60%)
- Highly Sensitive Model (FN = 0, TN @ 80%)
- Highly Specific Model (FP = 0, TP @ 20%)
- Highly Specific Model (FP = 0, TP @ 40%)
- Highly Specific Model (FP = 0, TP @ 60%)
- Highly Specific Model (FP = 0, TP @ 80%)
- Ideal Model (FP = FN = 0, @ 100%)
- Inverse Model (TP = TN = 0, -100%)
- Null Model (PP = FP = TP = 0)

** Incorporated cases with unbalanced error prediction rates*



**Flat lines reflect perfectly sensitive model*



Note many cases undefined with zero values

Metrics calculated for hypothetical cases, varying class imbalance

Modified MCC

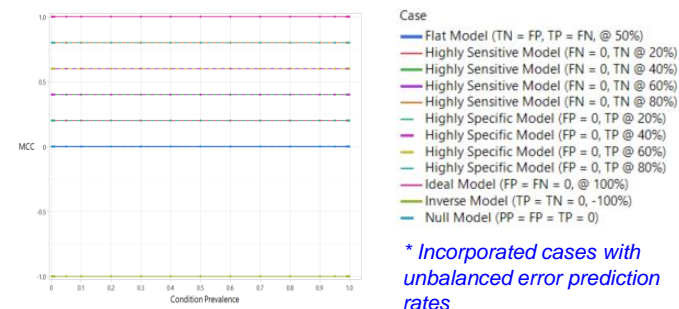
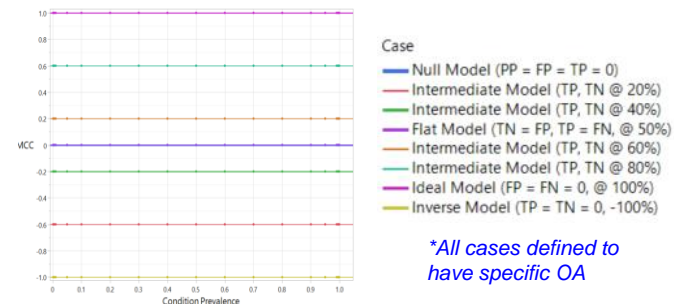


- Alternative calculation approaches can improve MCC
 - Scale every cell element to be out of 100%, calculate as a function of proportions
 - Or divide by $AP^2 \cdot AN^2$ instead of standard denominator

$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}$$

$$MCC^{*mod} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FN)^2 \cdot (TN + FP)^2}} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{AP^2 \cdot AN^2}}$$

Modified MCC/Delta Rates

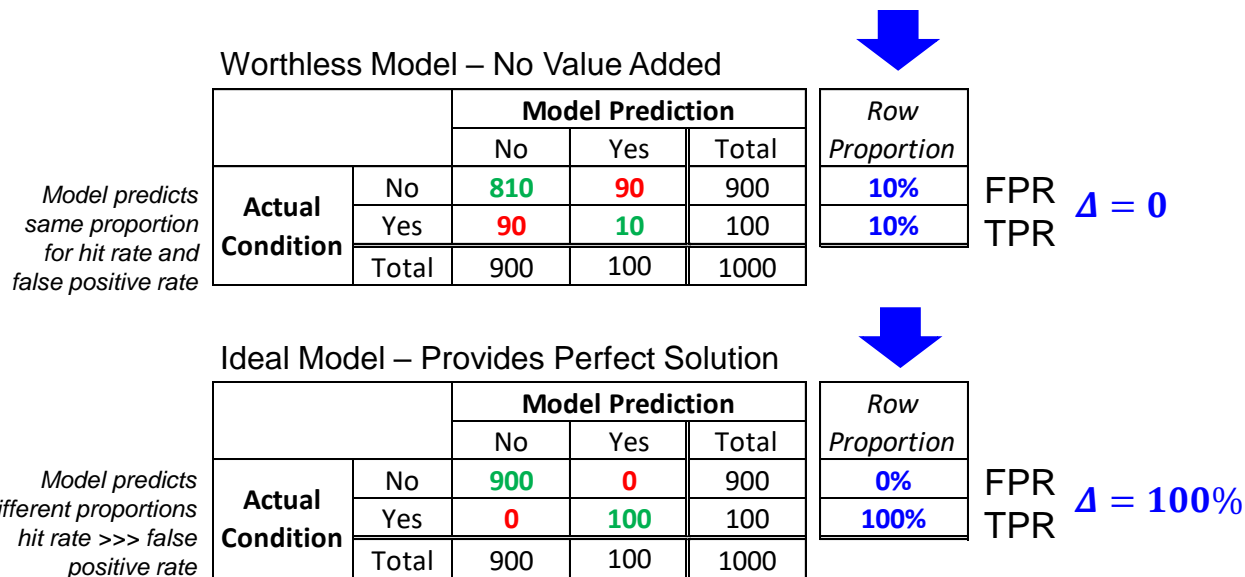


Adjustments stabilize MCC over full range of prevalence, equivalent to difference between True Positive to False Negative Rates

What is a “Good” Model?



- Adds value over the null model, provides some improved capability to differentiate between classes



FPR = False alarm rate, False Positive Rate

TPR = Recall, sensitivity, Probability of Detection (POD), hit rate

Focus on maximizing the difference between True & False Positive Rates

Importance can be assessed with standard statistical t-test of 2 proportions

Summary and Recommendations



- Many popular metrics used to assess classifier performance are affected by class imbalance, so impacts are widespread but unrecognized
- Be skeptical of claims made about classification model performance, especially if given in terms of “accuracy”
- A “good” model is one that adds value, provides some improved capability to differentiate between classes – and effects can be quantified
- Modifications can be made to stabilize metrics over full range of possible prevalence values, but model quality can be best quantified & assessed for statistical significance through the comparison of the true positive to true negative rates



Questions?