

# Decision Trees

**Practical Machine Learning (with R)**

UC Berkeley

Spring 2015



# Topics

- ⇒ Administrative
  - Role Call
  - Assignments due to github
  - Class Google Group (All joined)
- ⇒ Expectations (Review)
- ⇒ New Topics
  - R Meetup

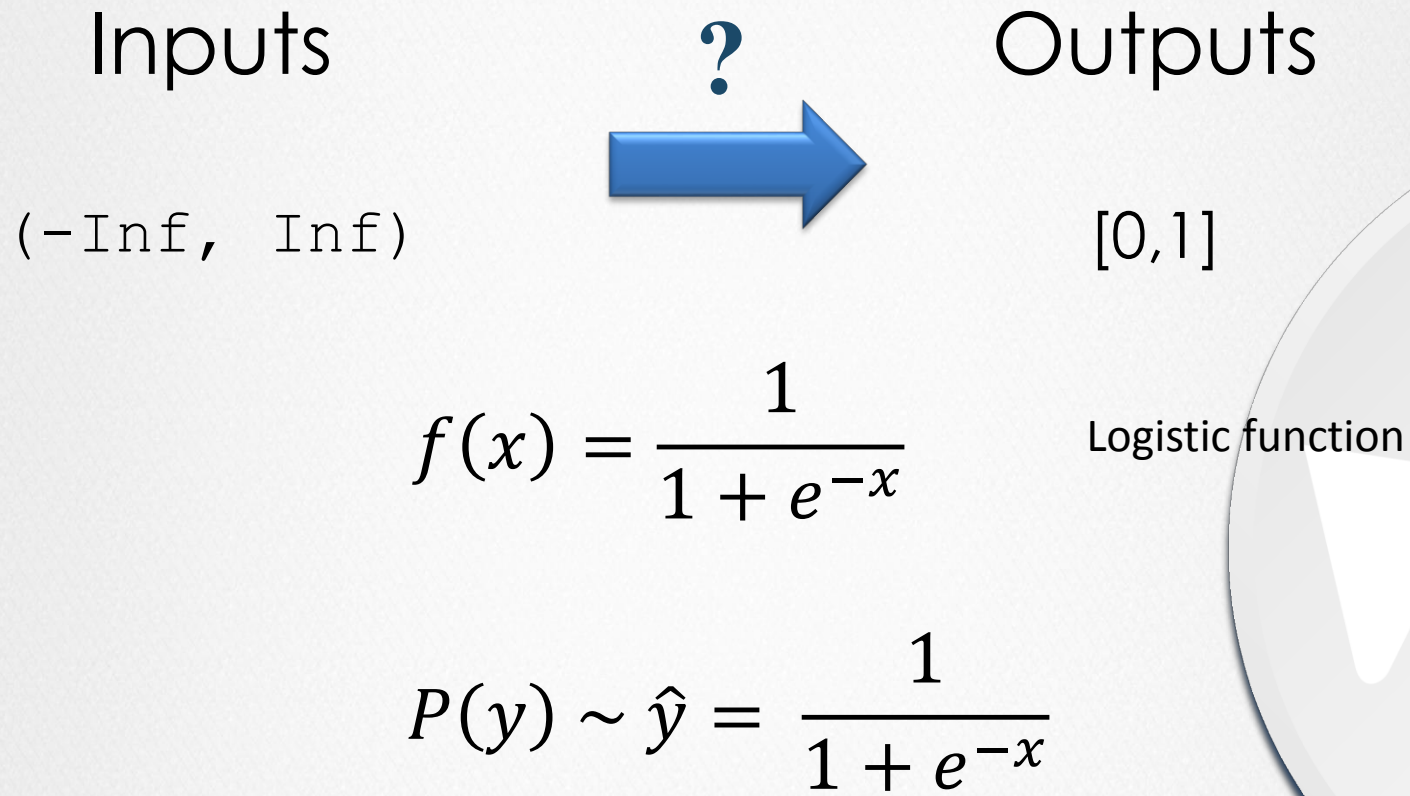


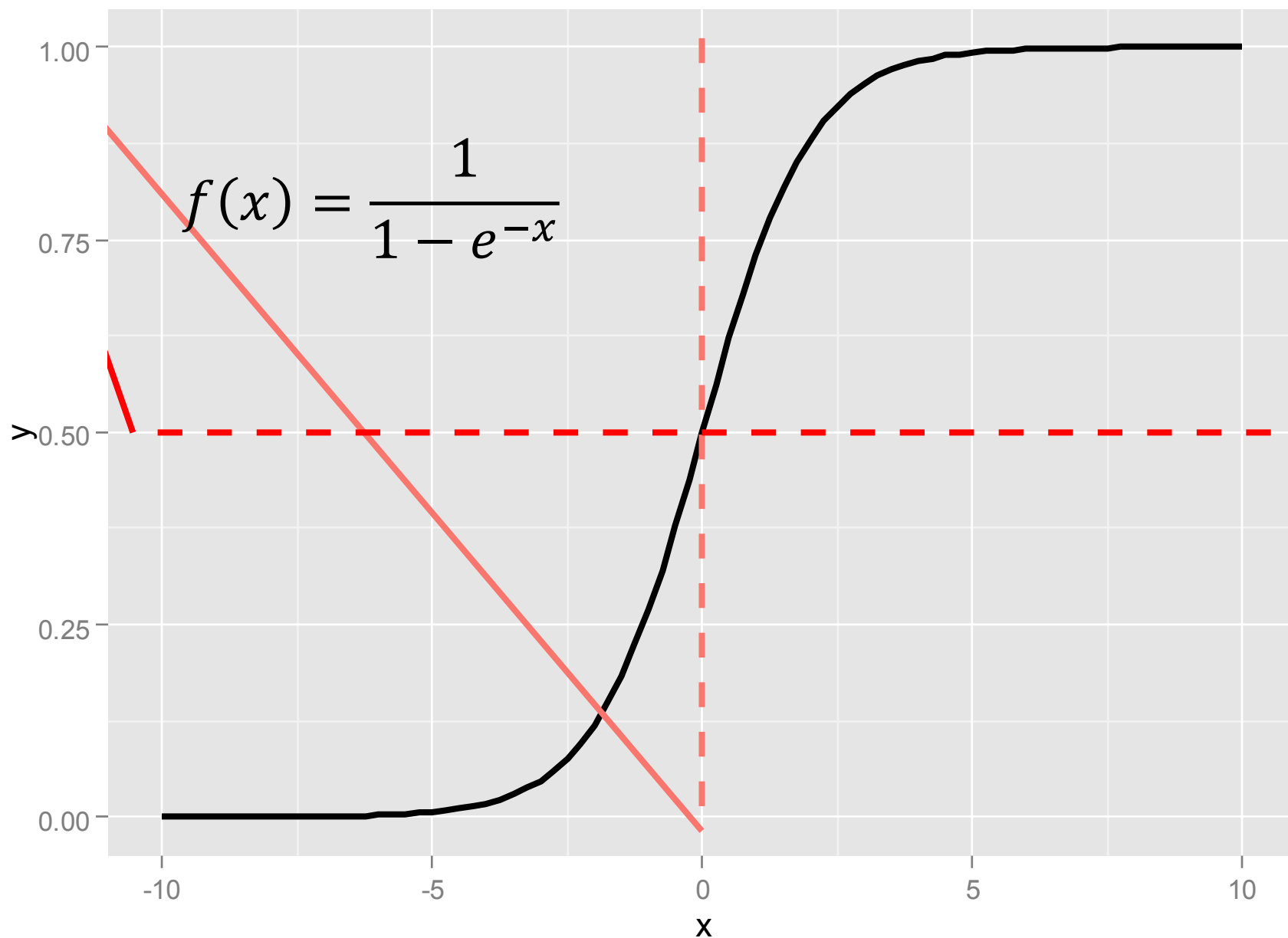
# REVIEW





# Need a tool ...





## Now WHAT

- ➔ Proceed as we would with linear regression ... and look for  $\beta$ 's

$$\hat{y} \sim \frac{1}{1 + e^{-x}}$$

$$\hat{y} \sim \frac{1}{1 + e^{-\beta_0 + \sum_{i=1}^p \beta_i x_i}}$$

- ➔ Then solve as linear regression:

$$\operatorname{argmin}_{\beta} \left( \sum (\hat{y} - y)^2 \right)$$



# LOGISTIC REGRESSION SUMMARY

Call:

```
glm(formula = Versicolor ~ . - Sepal.Length, family = binomial,  
     data = train)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.1262	-0.7731	-0.3984	0.8063	2.1562

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	6.9506	2.2261	3.122	0.00179	**
Sepal.Width	-2.9565	0.6668	-4.434	9.26e-06	***
Petal.Length	1.1252	0.4619	2.436	0.01484	*
Petal.Width	-2.6148	1.0815	-2.418	0.01562	*

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 190.95 on 149 degrees of freedom  
Residual deviance: 145.21 on 146 degrees of freedom  
AIC: 153.21

Number of Fisher Scoring iterations: 5

Log Odds

Variable

- Significance?
- Importance?





# MODEL PERFORMANCE





# Model Performance (thus far)

- ⇒ Determine performance metric:
  - **RMSE (regression)**
  - **Accuracy (classification)**
- ⇒ Fit Model
- ⇒ Calculate statistic (“metric”) on Training Data

“*training*” or “*apparent*” performance will:

- over-fit to training data
- predict very well, unbelievably well
- Not generalize to *new data*.



# CARDINAL RULE

**DO NOT ESTIMATE PERFORMANCE ON  
YOUR TRAINING DATA**

**→ Need tool for unbiased estimate for  
calculating performance**



# MEASUREMENTS AND STATISTICS

## Measurement

*Quantification of a phenomena*

## Statistic

*measurement of a stochastic phenomena*

Deterministic  
≠  
Stochastic

## Examples

- `mean(x)`      `<- x` is generated by a stochastic process
- `sd(x)`



# STATISTICS

- ⇒ “True” value unknown → uncertainty
- ⇒ Uncertainty can be measured
  - Variance
  - Standard deviation
  - Confidence Interval
  - ...
- ⇒ Repeated measurements decrease the uncertainty



## EXERCISE 1: CALCULATE `sd (mean (x) )`



# RESAMPLING

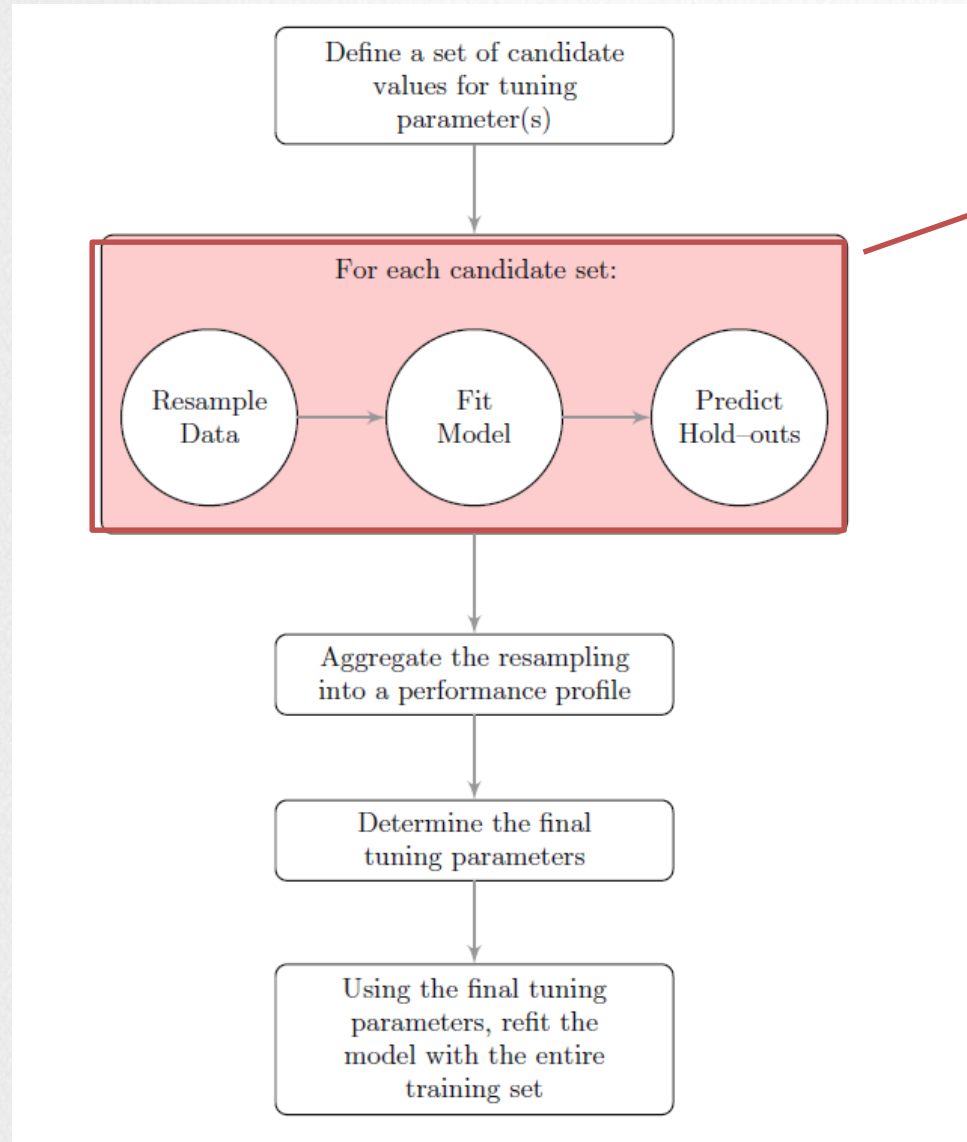
Kuhn benefits of resampling:

- Selection of optimal tuning parameter(s)  
“With so many choices how do we
- Unbiased estimate of model performance





# KUHN'S RESAMPLING PROCESS



Today's Focus



# RESAMPLING

- ➔ Best Solution (n-permitting)
  - split data into training and test data
  - and do what Kuhn says.

## Why(?)

- Easy to interpret defend
- Requires data not be consumed by model
- Computationally easy
- Is generally not (by itself) the most accurate → no confidence



# RESAMPLING STRATEGIES

- ➔ Repeated Splitting
- ➔ K-Fold Cross Validation
- ➔ Bootstrap





# REPEATED SPLITTING

## AKA Monte Carlo Splitting

- ⇒ Split data 75%-25%
  - Fit Model
  - Calculate Performance Metric
  - Repeat with Different Split (K-times)
- ⇒ Calculate Metric

$$Metric = AVG_i(metric)$$



# 10-Fold Cross Validation



...

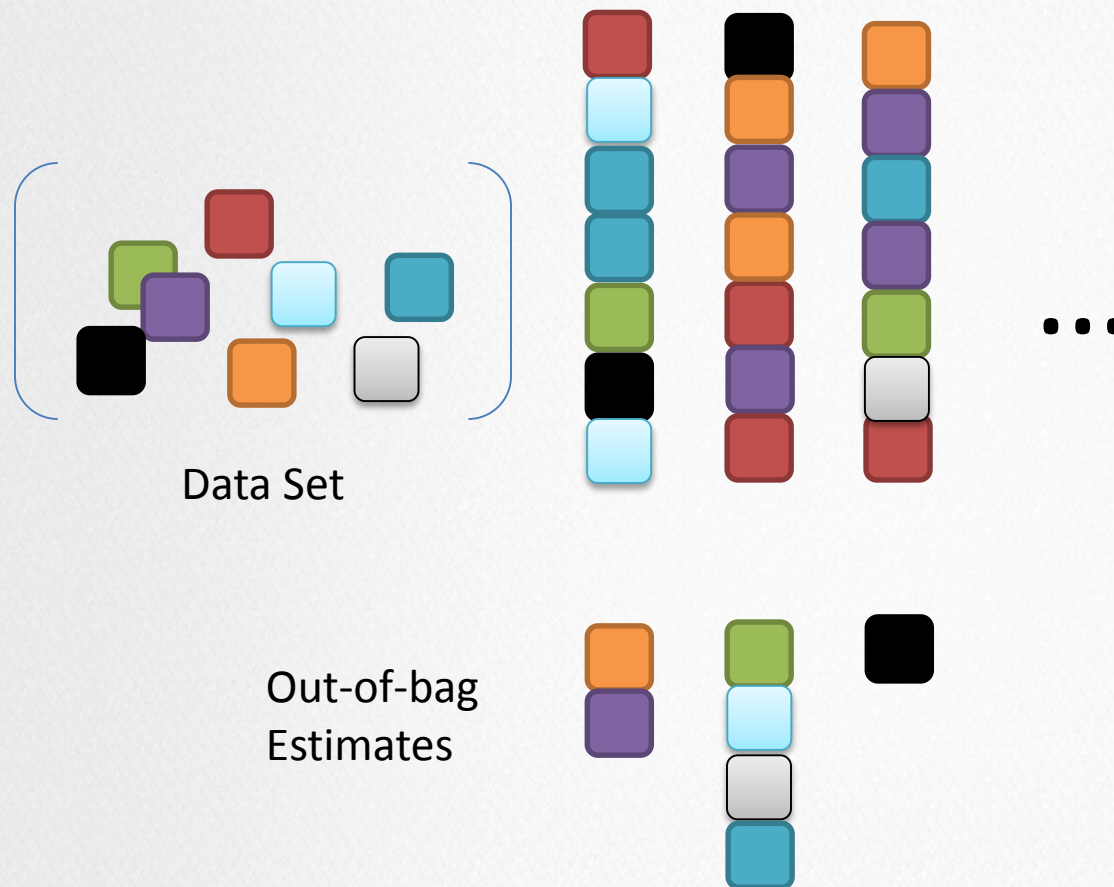


LOOCV :  $K \rightarrow n$

- Split the data set into 10 equal sized samples.
- Leave one sample out (fold)
  - Fit the model
  - calculate the metric on the fold
  - Repeat choosing another sample until done
- Calculate Metric
$$Metric = AVG_i(metric)$$
- 5 or 10-fold common

# Bootstrap

⇒ “Sampling with Replacement”





# Which Is Best?

→ There isn't one.

K-fold cross validation

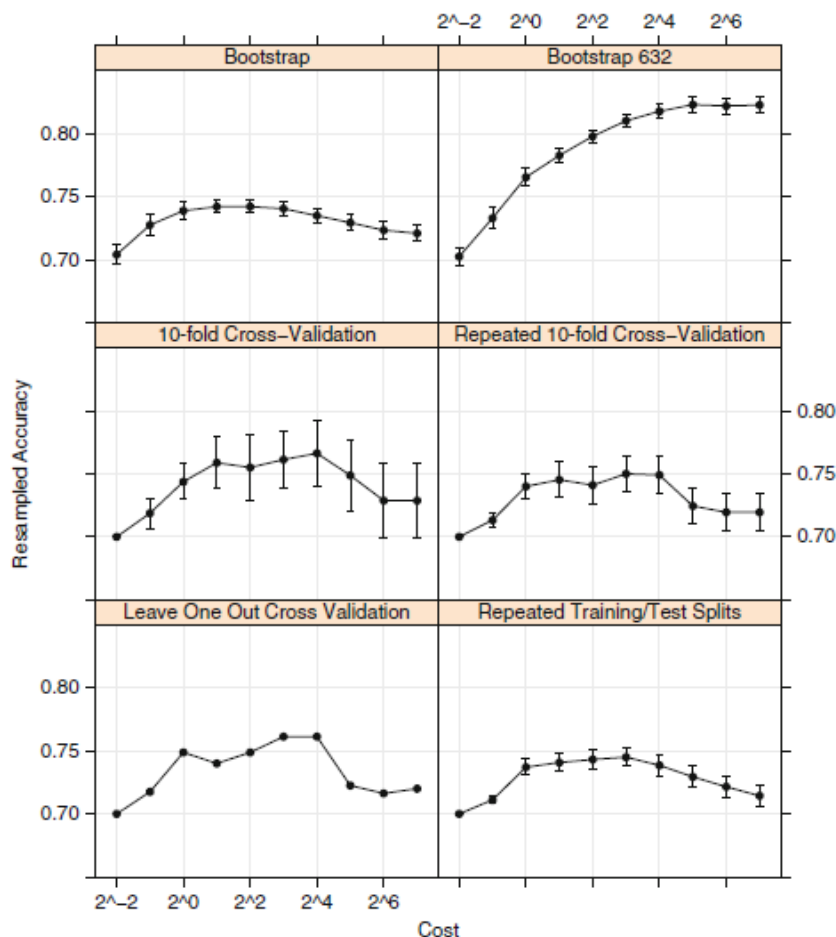
Higher Variance

Lower Bias

Bootstrap

Lower Variance

Higher Bias





**MODEL PERFORMANCE IS NOT  
TRAINING PERFORMANCE**



# CLASSIFICATION PERFORMANCE



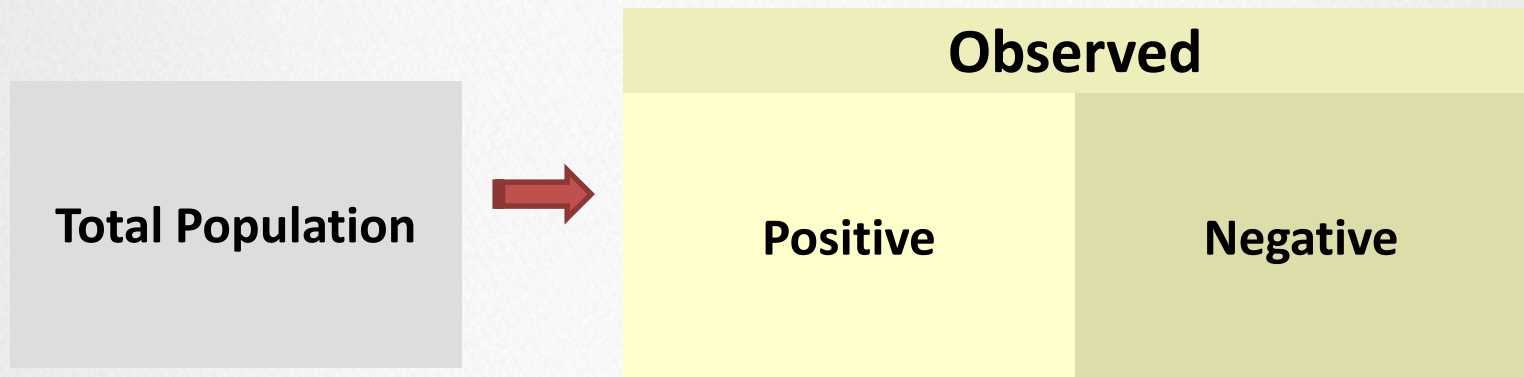


# METRICS FOR BI-NOMIAL CLASSIFICATION



**Total Population**







**Total Population**



**Predicted**

**Positive**

**Negative**



**Total Population**



<b>Observed</b>	
<b>Positive</b>	<b>Negative</b>



<b>Predicted</b>	<b>Positive</b>
	<b>Negative</b>

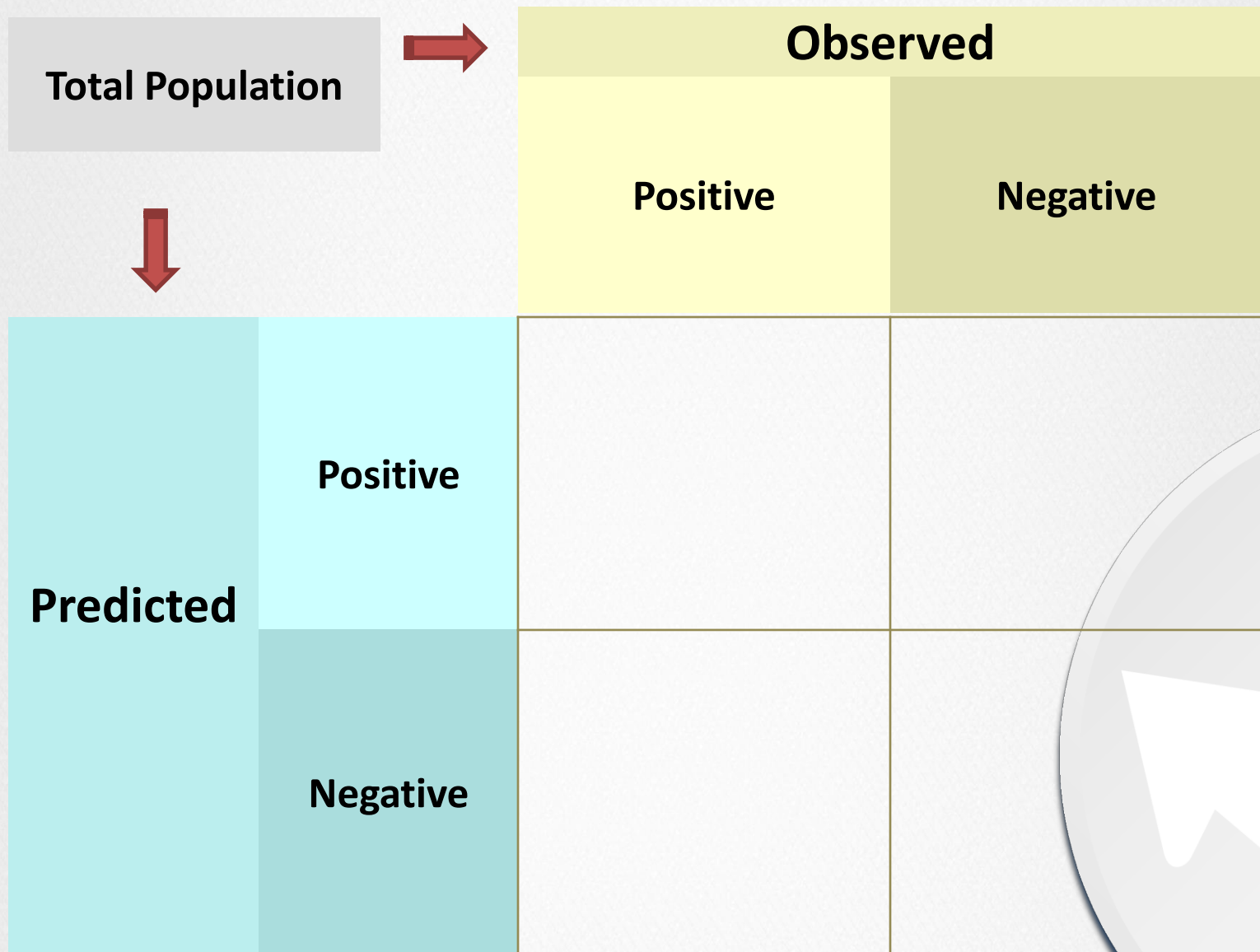
## Accuracy

“How many did I get right”

## Error Rate

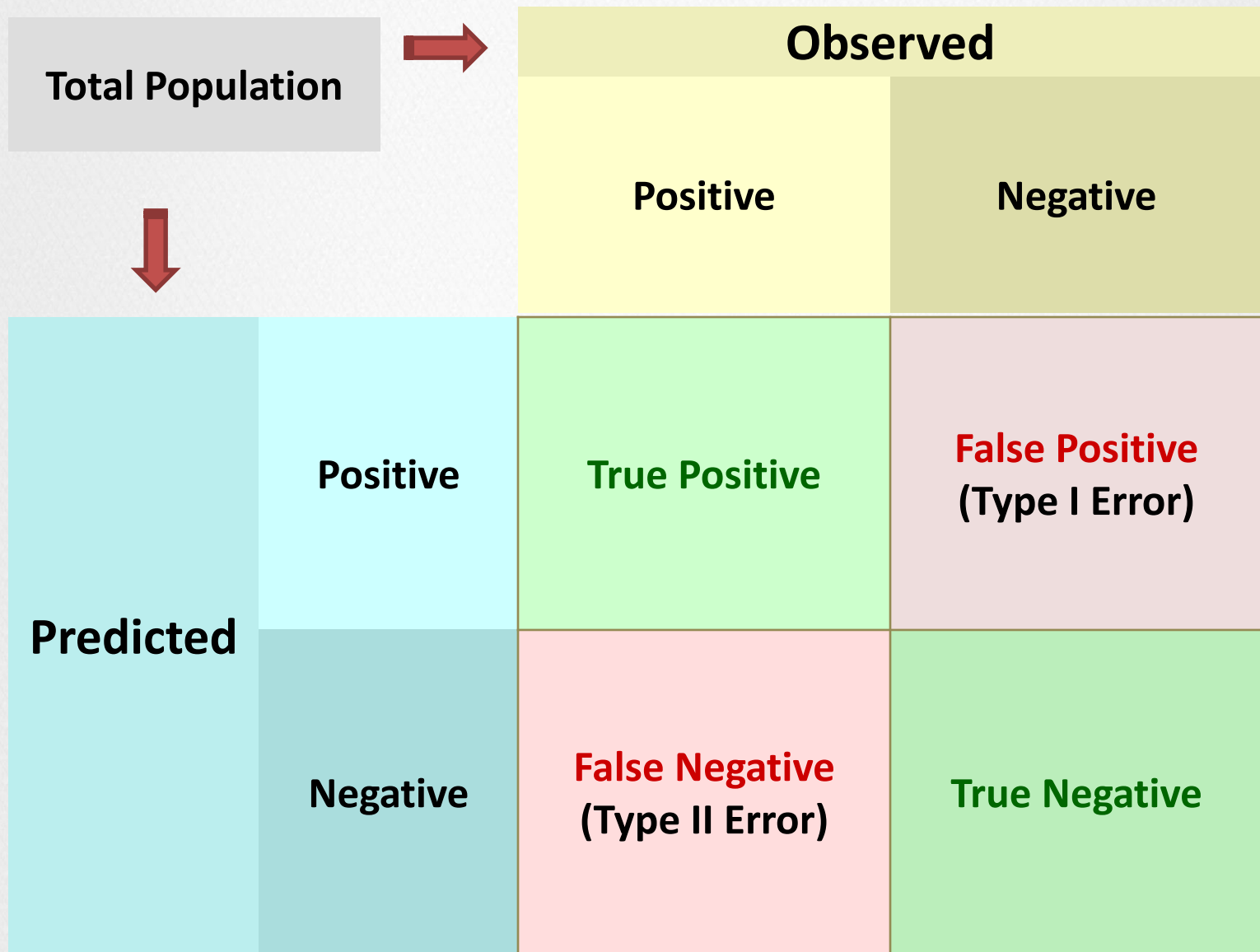
or Misclassification Rate  
“How many did I get wrong”





- [https://en.wikipedia.org/wiki/Sensitivity\\_and\\_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)





# Alternatives: Norm by Observed

		Observed	
		Positive	Negative
Predicted	Positive	True Positive Rate (TPR), <b>Sensitivity</b> , Recall  $\frac{\text{True Positives}}{\text{Observed Positives}}$	False Positive Rate (FPR), Fall-Out  $\frac{\text{False Positives}}{\text{Observed Negatives}}$
	Negative	False Neg. Rate (FNR), Miss rate  $\frac{\text{False Negatives}}{\text{Observed Positives}}$	True Neg. Rate (TNR), <b>Specificity</b> (SPC)  $\frac{\text{True Negatives}}{\text{Observed Negatives}}$

# Alternatives: Norm by Predicted

		Observed	
		Positive	Negative
Predicted	Positive	Pos. Predictive Value (PPV), <b>Precision</b> $\frac{\text{True Positives}}{\text{Predicted Positives}}$	False Discovery Rate (FDR) $\frac{\text{False Positives}}{\text{Predicted Positives}}$
	Negative	False Omission Rate (FOR) $\frac{\text{False Negatives}}{\text{Predicted Negatives}}$	Negative Predictive Value (NPV) $\frac{\text{True Negatives}}{\text{Predicted Negatives}}$

- [https://en.wikipedia.org/wiki/Sensitivity\\_and\\_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)



# MORE FUN ...

[https://en.wikipedia.org/wiki/Sensitivity\\_and\\_specificity](https://en.wikipedia.org/wiki/Sensitivity_and_specificity)



## EXERCISE: BINOMIAL METRICS



## EVEN MORE COMPLICATION

- ➡ Not all errors need count “equivocal zone” or “intermediate zone”
- ➡ *Prevalent when the model has three choices, e.g. A or B or Nothing.*





# MUTLINOMIAL CLASSIFICATION



# TERMS

- ⇒ Kappa Statistic,
  - ⇒ S-Statistics, F-Statistic
- 



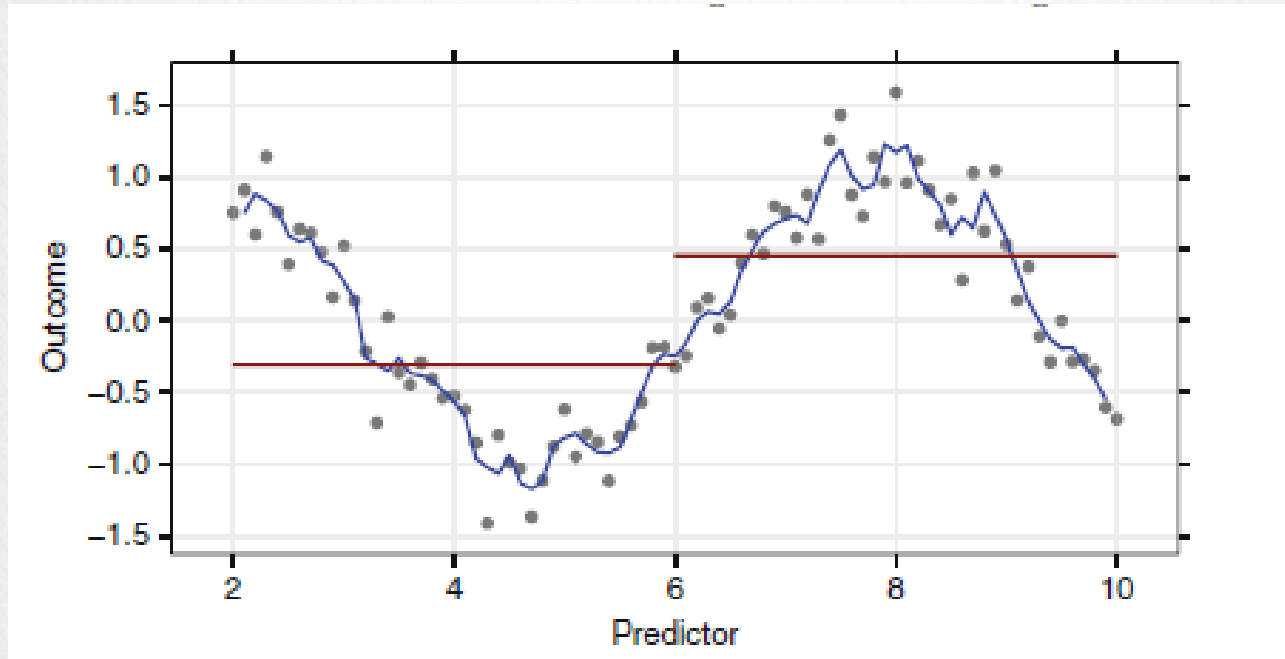
# MULTICLASS CLASSIFICATION WITH LOGISTIC REGRESSION





# BIAS VARIANCE TRADE-OFF

$$E[MSE] = \sigma^2 + (\text{model bias})^2 + \text{model variance}$$



		True condition			
		Condition positive	Condition negative	$\text{Prevalence} = \frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	
Predicted condition	Predicted condition positive	<u>True positive</u>	<u>False positive</u> (Type I error)	<u>Positive predictive value</u> (PPV), <u>Precision</u> = $\frac{\sum \text{True positive}}{\sum \text{Test outcome positive}}$	<u>False discovery rate</u> (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Test outcome positive}}$
	Predicted condition negative	<u>False negative</u> (Type II error)	<u>True negative</u>	<u>False omission rate</u> (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Test outcome negative}}$	<u>Negative predictive value</u> (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Test outcome negative}}$
		<u>True positive rate</u> (TPR), <u>Sensitivity</u> , Recall = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	<u>False positive rate</u> (FPR), <u>Fall-out</u> = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	<u>Positive likelihood ratio</u> (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	
		<u>Accuracy</u> (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$			<u>Diagnostic odds ratio</u> (DOR) = $\frac{\text{LR+}}{\text{LR-}}$
		<u>False negative rate</u> (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	<u>True negative rate</u> (TNR), <u>Specificity</u> (SPC) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	<u>Negative likelihood ratio</u> (LR-) = $\frac{\text{FNR}}{\text{TNR}}$	