

# Generalizability and Data Splitting

---

February 24, 2026

# Mid-Semester Feedback Survey

---

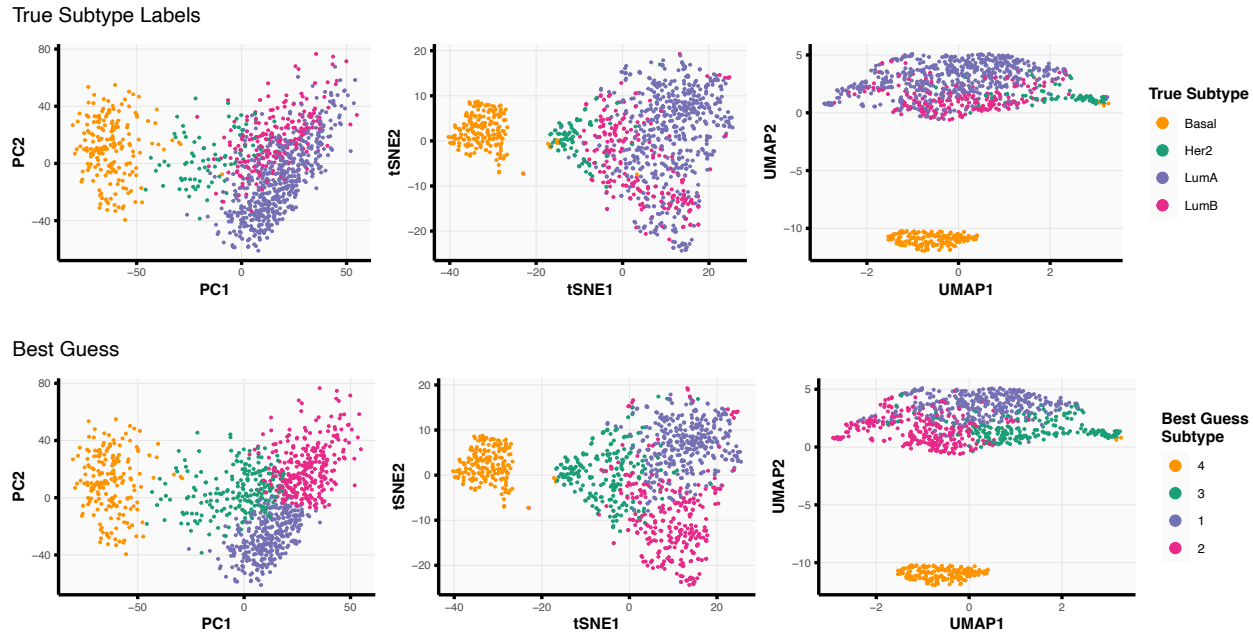
Please fill out anonymous mid-semester feedback survey:

<https://forms.gle/KoXjt3CbY91GKk7V8>



# Lab 2 Recap

- + Lab 2 solutions have been released along with true breast cancer subtypes
  - See your `dsip/lab2/results/true_clusters_plot.pdf` for a quick visualization



# Today's plan: Introduction to Supervised Learning

---

- 1 **Review** of last time
- 2 **Evaluation metrics** to assess generalizability
- 3 **Data splitting**
- 3 **Introduction to Lab 3**

## Last Time: **The #1 Goal** of Supervised Learning

The goal of supervised learning is to learn insights from the current data that are **generalizable to future, unseen data**

**Generalizability** is closely related to:

- + **Bias of the model:** on average, how wrong is the model compared to the truth?
- + **Variance of the model:** if I retrained my model on a similar but new dataset, how much does my model change?

$$\underbrace{\mathbb{E}[(y - \hat{f}(x))^2]}_{\substack{\text{population MSE} \\ \text{(generalization error)}}} = \left( \text{Bias}(\hat{f}(x)) \right)^2 + \text{Var}(\hat{f}(x)) + \underbrace{\sigma^2}_{\substack{\text{irreproducible} \\ \text{error}}}$$

# Regression and Classification Methods

|               | Linear/<br>Logistic   | Ridge   | LASSO   | Elastic Net   | Decision<br>Trees  | Random<br>Forests  | Neural<br>Networks   |
|---------------|---|---|---|---|--|--|--|
| Advantages    | <ul style="list-style-type: none"><li>• Simple</li><li>• Can do inference (with the usual assumptions)</li></ul>      | <ul style="list-style-type: none"><li>• Good with correlated features</li><li>• MSE Existence Theorem</li></ul>             | <ul style="list-style-type: none"><li>• Feature selection and sparsity</li><li>• Interpretable</li></ul>  | <ul style="list-style-type: none"><li>• Compromise between ridge and LASSO</li></ul>                | <ul style="list-style-type: none"><li>• Interpretable</li><li>• Non-linear</li></ul>             | <ul style="list-style-type: none"><li>• Non-linear</li><li>• Powerful predictors</li></ul>   | <ul style="list-style-type: none"><li>• Can naturally model very complex data structures</li><li>• Powerful predictors</li></ul> |
| Disadvantages | <ul style="list-style-type: none"><li>• Can be unstable</li><li>• Major problems when <math>p &gt; n</math></li></ul> | <ul style="list-style-type: none"><li>• Dense model is not interpretable</li><li>• No automatic feature selection</li></ul> | <ul style="list-style-type: none"><li>• Often selects 1 feature from correlated group</li><li>• Gives lower accuracy if truth is not sparse</li></ul> | <ul style="list-style-type: none"><li>• More difficult to tune two parameters in practice</li></ul> | <ul style="list-style-type: none"><li>• Unstable</li><li>• Greedy, axis-aligned splits</li></ul> | <ul style="list-style-type: none"><li>• Not as interpretable as decision trees</li><li>• Not efficient at learning smooth structures</li></ul> | <ul style="list-style-type: none"><li>• Difficult to tune</li><li>• Generally requires lots of data</li></ul>                    |

and more... (e.g., K nearest neighbors, support vector machines, bagging, ensembling)

## Today: Revisiting **the #1 Goal** of Supervised Learning

---

The goal of supervised learning is to learn insights from the current data that are **generalizable to future, unseen data**

**Q:** How do we accurately **assess** generalizability (i.e., how well our model will do on new data)?

- + Prediction performance metrics
- + Data splitting

# Prediction Performance Metrics

---

# Regression Metrics

---

**True (observed)** responses:  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

**Predicted** responses:  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$

---

| Metric | Definition | When to use |
|--------|------------|-------------|
|--------|------------|-------------|

---

# Regression Metrics

---

**True (observed) responses:**  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

**Predicted responses:**  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$

---

| Metric                          | Definition                                     | When to use   |
|---------------------------------|--|---|
| <b>Mean squared error (MSE)</b> | $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ | - Most popular and arguably the default metric for regression |

# Regression Metrics

**True (observed) responses:**  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

**Predicted responses:**  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$

| Metric                                | Definition  | When to use   |
|---------------------------------------|---|---|
| <b>Mean squared error (MSE)</b>       | $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$        | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li></ul>   |
| <b>Root mean squared error (RMSE)</b> | $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$ | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li><li>- Error measured on same scale/units as y</li></ul> |

# Regression Metrics

**True (observed) responses:**  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

**Predicted responses:**  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$

| Metric                                | Definition  | When to use   |
|---------------------------------------|---|---|
| <b>Mean squared error (MSE)</b>       | $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li></ul>   |
| <b>Root mean squared error (RMSE)</b> | $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$   | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li><li>- Error measured on same scale/units as y</li></ul> |
| <b>R<sup>2</sup></b>                  | $1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{n}{n-1} \frac{MSE(\mathbf{y}, \hat{\mathbf{y}})}{\hat{Var}(\mathbf{y})}$ | <ul style="list-style-type: none"><li>- Normalized MSE (normalized by variance of observed responses)</li></ul>   |

# Regression Metrics

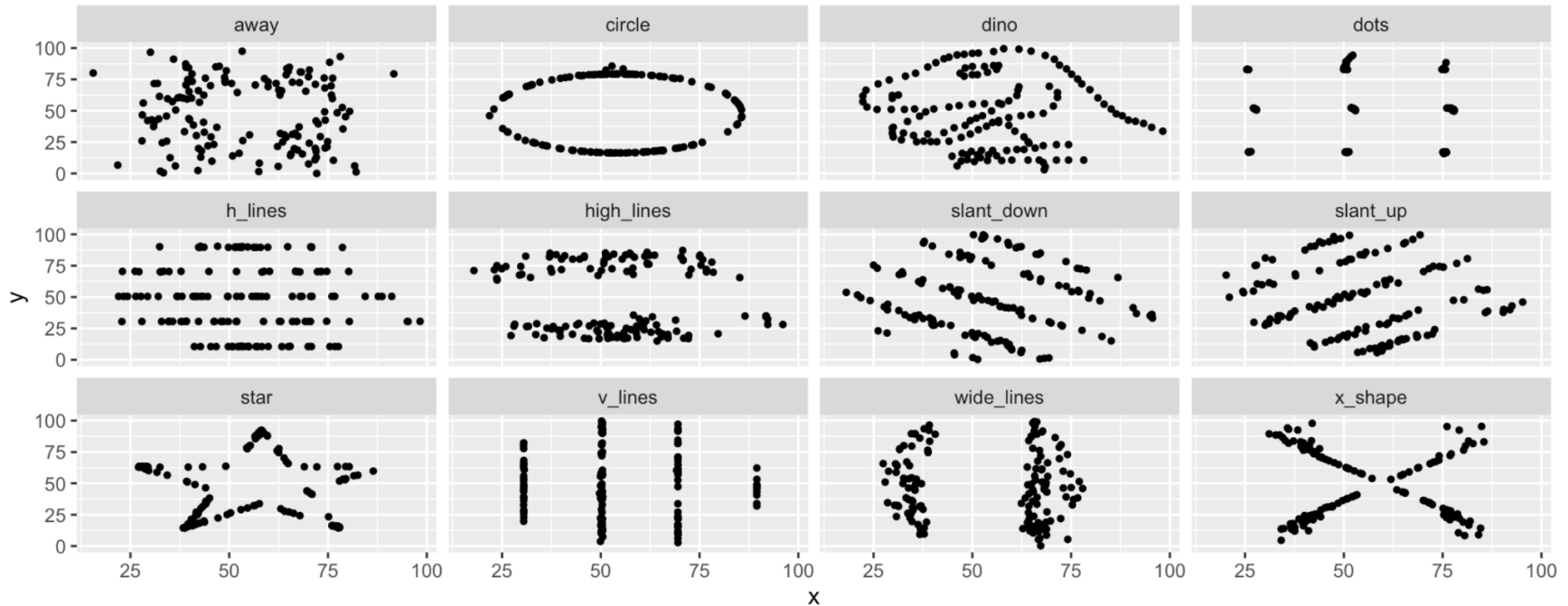
**True (observed) responses:**  $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$

**Predicted responses:**  $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_n)^\top \in \mathbb{R}^n$

| Metric                                | Definition  | When to use   |
|---------------------------------------|---|---|
| <b>Mean squared error (MSE)</b>       | $\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li></ul>   |
| <b>Root mean squared error (RMSE)</b> | $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$   | <ul style="list-style-type: none"><li>- Most popular and arguably the default metric for regression</li><li>- Error measured on same scale/units as y</li></ul> |
| <b>R<sup>2</sup></b>                  | $1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{n}{n-1} \frac{MSE(\mathbf{y}, \hat{\mathbf{y}})}{\hat{Var}(\mathbf{y})}$ | <ul style="list-style-type: none"><li>- Normalized MSE (normalized by variance of observed responses)</li></ul>   |
| <b>Mean absolute error (MAE)</b>      | $\frac{1}{n} \sum_{i=1}^n  y_i - \hat{y}_i $  | <ul style="list-style-type: none"><li>- More robust to outliers than MSE/RMSE</li></ul>   |

# We need to go beyond a single summary statistic

Each of these has the same mean, standard deviation, variance, and correlation

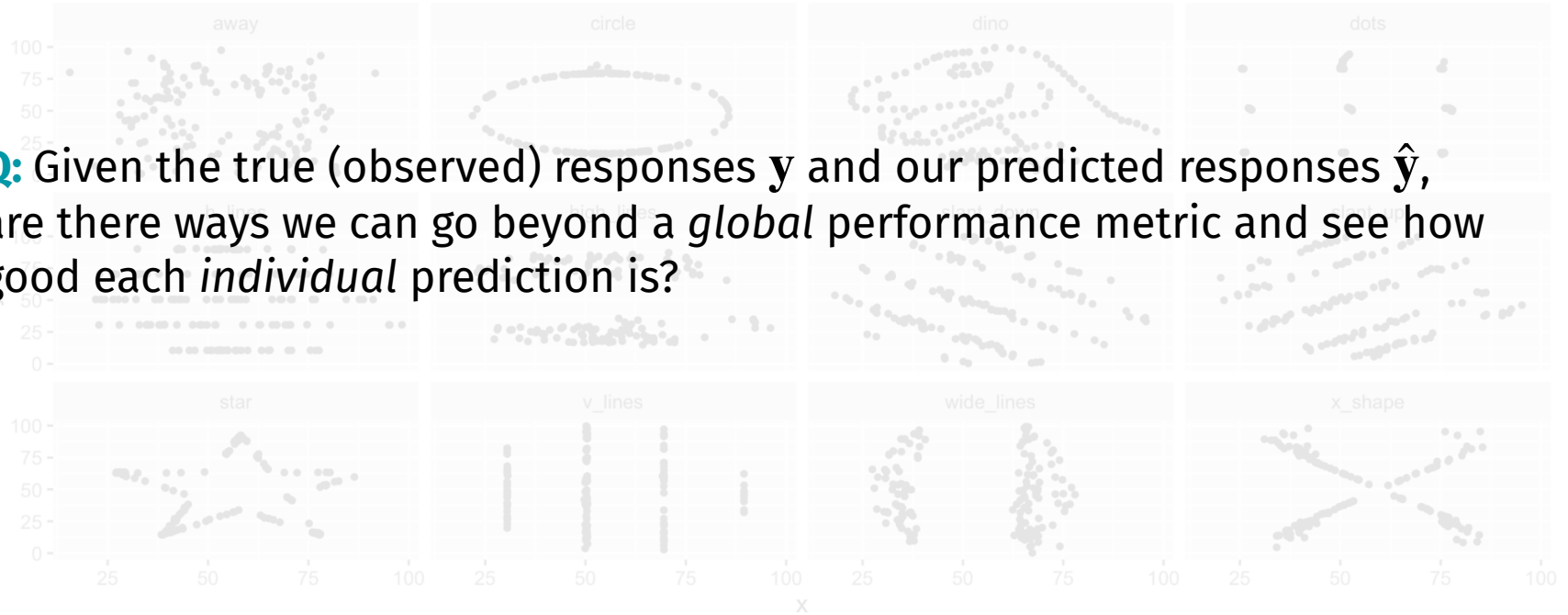


"The Datasaurus Dozen" [[Matejka and Fitzmaurice \(2017\)](#)]

# We need to go beyond a single summary statistic

Each of these has the same mean, standard deviation, variance, and correlation

**Q:** Given the true (observed) responses  $\mathbf{y}$  and our predicted responses  $\hat{\mathbf{y}}$ , are there ways we can go beyond a *global* performance metric and see how good each *individual* prediction is?

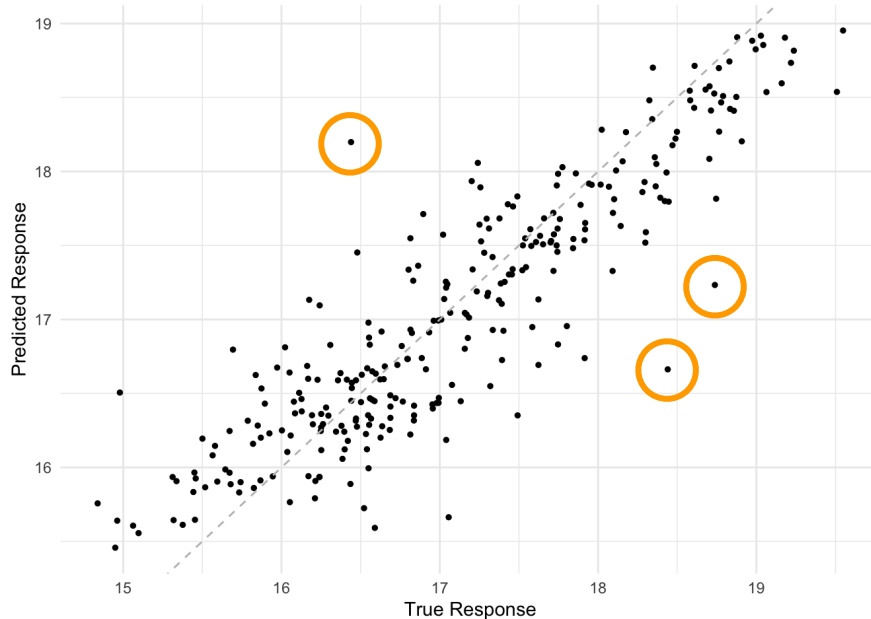


"The Datasaurus Dozen" [\[Matejka and Fitzmaurice \(2017\)\]](#)

# Going beyond "global" summary metrics...

Plot true (observed) responses  $\mathbf{y}$  versus predicted responses  $\hat{\mathbf{y}}$

- + Allows for further investigation into "odd" observations



# Output of classification methods

---

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

# Output of classification methods

| True (observed)<br>class labels                                   | Predicted<br>class labels   | Predicted<br>class probabilities  |
|---|---|---|
| $\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$ | $\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$ | $\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$ |

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$\hat{y}_i \in \{c_1, \dots, c_K\}$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$\hat{y}_i \in \{c_1, \dots, c_K\}$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

↑  
Probability of being  
in class K

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

↑  
Probability of being  
in class K

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = 1) \\ \vdots \\ Pr(\hat{y}_n = 1) \end{bmatrix}$$

For binary classification, this is typically rewritten as a column vector (not an  $n \times 2$  matrix)

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

↑  
Probability of being  
in class K

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = 1) \\ \vdots \\ Pr(\hat{y}_n = 1) \end{bmatrix}$$

For binary classification, this is typically rewritten as a column vector (not an  $n \times 2$  matrix)

# Classification Accuracy

---

$$\begin{aligned} + \text{ Classification accuracy} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i = \hat{y}_i\} \\ &= \frac{\# \text{ of correctly predicted class labels}}{\text{total } \# \text{ of samples}} \end{aligned}$$

# Classification Accuracy

---

+ **Classification accuracy**  $= \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i = \hat{y}_i\}$   
 $= \frac{\# \text{ of correctly predicted class labels}}{\text{total } \# \text{ of samples}}$

+ Higher is better

# Classification Accuracy

---

$$\begin{aligned} + \text{ Classification accuracy} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i = \hat{y}_i\} \\ &= \frac{\# \text{ of correctly predicted class labels}}{\text{total } \# \text{ of samples}} \end{aligned}$$

- + Higher is better
- + What if we have class imbalance?

# Classification Accuracy

---

$$\begin{aligned} + \text{ Classification accuracy} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{y_i = \hat{y}_i\} \\ &= \frac{\# \text{ of correctly predicted class labels}}{\text{total } \# \text{ of samples}} \end{aligned}$$

+ Higher is better

+ What if we have class imbalance?

- Ex. Suppose we have a sample of 100 people and only 10 have the disease. If we always predict healthy, we get 90% classification accuracy!

# Confusion Matrix

We can go beyond classification accuracy and break the predicted responses down into a **confusion matrix**

|                  |                             | Predicted condition                                 |   |
|------------------|-----------------------------|---|---|
|                  |                             | Predicted positive (PP)                             | Predicted negative (PN)                                 |
| Actual condition | Positive (P)<br>[a]         | True positive (TP),<br>hit <sup>[b]</sup>           | False negative (FN),<br>miss,<br>underestimation        |
|                  | Negative (N) <sup>[d]</sup> | False positive (FP),<br>false alarm, overestimation | True negative (TN),<br>correct rejection <sup>[e]</sup> |

# Other classification metrics based on **predicted classes**

Sources: [16][17][18][19][20][21][22][23] view · talk · edit

|   |  | Predicted condition   |   |   |   |
|---|--|---|---|---|---|
|   |  | Predicted positive (PP)   | Predicted negative (PN)   |   |   |
| Actual condition  | Total population<br>= P + N  |   |   | Informedness, bookmaker informedness (BM)<br>= TPR + TNR - 1  | Prevalence threshold (PT)<br>$= \frac{(\text{PT})}{\sqrt{\text{TPR} \times \text{FPR}} - \text{FPR}} - \text{FPR}$    |
|   | Positive (P)<br>[a]  | True positive (TP), hit <sup>[b]</sup>  | False negative (FN), miss, underestimation  | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power<br>$= \frac{\text{TP}}{\text{P}} = 1 - \text{FNR}$ | False negative rate (FNR), miss rate, type II error <sup>[c]</sup><br>$= \frac{\text{FN}}{\text{P}} = 1 - \text{TPR}$ |
|   | Negative (N) <sup>[d]</sup>  | False positive (FP), false alarm, overestimation                                    | True negative (TN), correct rejection <sup>[e]</sup>  | False positive rate (FPR), probability of false alarm, fall-out, type I error <sup>[f]</sup><br>$= \frac{\text{FP}}{\text{N}} = 1 - \text{TNR}$   | True negative rate (TNR), specificity (SPC), selectivity<br>$= \frac{\text{TN}}{\text{N}} = 1 - \text{FPR}$           |
| Prevalence<br>$= \frac{\text{P}}{\text{P} + \text{N}}$                  | Positive predictive value (PPV), precision<br>$= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$   | False omission rate (FOR)<br>$= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$       | Positive likelihood ratio (LR+)<br>$= \frac{\text{TPR}}{\text{FPR}}$  | Negative likelihood ratio (LR-)<br>$= \frac{\text{FNR}}{\text{TNR}}$  |   |
| Accuracy (ACC)<br>$= \frac{\text{TP} + \text{TN}}{\text{P} + \text{N}}$ | False discovery rate (FDR)<br>$= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$   | Negative predictive value (NPV)<br>$= \frac{\text{TN}}{\text{PN}} = 1 - \text{FOR}$ | Markedness (MK), deltaP ( $\Delta p$ )<br>= PPV + NPV - 1   | Diagnostic odds ratio (DOR)<br>$= \frac{\text{LR}+}{\text{LR}-}$  |   |
| Balanced accuracy (BA)<br>$= \frac{\text{TPR} + \text{TNR}}{2}$         | F <sub>1</sub> score<br>$= \frac{2 \text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}} = \frac{2 \text{TP}}{2 \text{TP} + \text{FP} + \text{FN}}$ | Fowlkes–Mallows index (FM)<br>$= \sqrt{\text{PPV} \times \text{TPR}}$               | Matthews correlation coefficient (MCC)<br>$= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}}}{\sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}$ | Threat score (TS), critical success index (CSI), Jaccard index<br>$= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$                         |   |

# Other classification metrics based on predicted classes

Sources: [16][17][18][19][20][21][22][23] view · talk · edit

|                  |   | Predicted condition   |  |   |  |
|------------------|---|---|--|---|--|
|                  |   | Predicted positive (PP)   | Predicted negative (PN)  |   |  |
| Actual condition | Total population<br>= P + N                       |   |  | Informedness, bookmaker informedness (BM)<br>= TPR + TNR - 1  | Prevalence threshold (PT)<br>$= \frac{PT}{\sqrt{TPR \times FPR} - FPR}$                          |
|                  | Positive (P)<br>[a]                               | True positive (TP), hit <sup>[b]</sup>  | False negative (FN), miss, underestimation                     | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power<br>$= \frac{TP}{P} = 1 - FNR$                  | False negative rate (FNR), miss rate, type II error <sup>[c]</sup><br>$= \frac{FN}{P} = 1 - TPR$ |
|                  | Negative (N) <sup>[d]</sup>                       | False positive (FP), false alarm, overestimation  | True negative (TN), correct rejection <sup>[e]</sup>           | False positive rate (FPR), probability of false alarm, fall-out, type I error <sup>[f]</sup><br>$= \frac{FP}{N} = 1 - TNR$                    | True negative rate (TNR), specificity (SPC), selectivity<br>$= \frac{TN}{N} = 1 - FPR$           |
|                  | Prevalence<br>$= \frac{P}{P + N}$                 | Positive predictive value (PPV), precision<br>$= \frac{TP}{PP} = 1 - FDR$                         | False omission rate (FOR)<br>$= \frac{FN}{PN} = 1 - NPV$       | Positive likelihood ratio (LR+)<br>$= \frac{TPR}{FPR}$  | Negative likelihood ratio (LR-)<br>$= \frac{FNR}{TNR}$   |
|                  | Accuracy (ACC)<br>$= \frac{TP + TN}{P + N}$       | False discovery rate (FDR)<br>$= \frac{FP}{PP} = 1 - PPV$   | Negative predictive value (NPV)<br>$= \frac{TN}{PN} = 1 - FOR$ | Markedness (MK), deltaP ( $\Delta p$ )<br>= PPV + NPV - 1   | Diagnostic odds ratio (DOR)<br>$= \frac{LR+}{LR-}$   |
|                  | Balanced accuracy (BA)<br>$= \frac{TPR + TNR}{2}$ | F <sub>1</sub> score<br>$= \frac{2 PPV \times TPR}{PPV + TPR}$<br>$= \frac{2 TP}{2 TP + FP + FN}$ | Fowlkes–Mallows index (FM)<br>$= \sqrt{PPV \times TPR}$        | Matthews correlation coefficient (MCC)<br>$= \frac{\sqrt{TPR \times TNR \times PPV \times NPV}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$ | Threat score (TS), critical success index (CSI), Jaccard index<br>$= \frac{TP}{TP + FN + FP}$    |

# Other classification metrics based on **predicted classes**

Sources: [16][17][18][19][20][21][22][23] view · talk · edit

|                  |   | Predicted condition   |  |   |  |
|------------------|---|---|--|---|--|
|                  |   | Predicted positive (PP)   | Predicted negative (PN)  |   |  |
| Actual condition | Total population<br>= P + N                       |   |  | Informedness, bookmaker informedness (BM)<br>= TPR + TNR - 1  | Prevalence threshold (PT)<br>$= \frac{TPR \times FPR - FPR}{TPR - FPR}$                          |
|                  | Positive (P)<br>[a]                               | True positive (TP), hit <sup>[b]</sup>  | False negative (FN), miss, underestimation                     | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power<br>$= \frac{TP}{P} = 1 - FNR$                  | False negative rate (FNR), miss rate, type II error <sup>[c]</sup><br>$= \frac{FN}{P} = 1 - TPR$ |
|                  | Negative (N) <sup>[d]</sup>                       | False positive (FP), false alarm, overestimation  | True negative (TN), correct rejection <sup>[e]</sup>           | False positive rate (FPR), probability of false alarm, fall-out, type I error <sup>[f]</sup><br>$= \frac{FP}{N} = 1 - TNR$                    | True negative rate (TNR), specificity (SPC), selectivity<br>$= \frac{TN}{N} = 1 - FPR$           |
|                  | Prevalence<br>$= \frac{P}{P + N}$                 | Positive predictive value (PPV), precision<br>$= \frac{TP}{PP} = 1 - FDR$                         | False omission rate (FOR)<br>$= \frac{FN}{PN} = 1 - NPV$       | Positive likelihood ratio (LR+)<br>$= \frac{TPR}{FPR}$  | Negative likelihood ratio (LR-)<br>$= \frac{FNR}{TNR}$   |
|                  | Accuracy (ACC)<br>$= \frac{TP + TN}{P + N}$       | False discovery rate (FDR)<br>$= \frac{FP}{PP} = 1 - PPV$   | Negative predictive value (NPV)<br>$= \frac{TN}{PN} = 1 - FOR$ | Markedness (MK), deltaP ( $\Delta p$ )<br>= PPV + NPV - 1   | Diagnostic odds ratio (DOR)<br>$= \frac{LR+}{LR-}$   |
|                  | Balanced accuracy (BA)<br>$= \frac{TPR + TNR}{2}$ | F <sub>1</sub> score<br>$= \frac{2 PPV \times TPR}{PPV + TPR}$<br>$= \frac{2 TP}{2 TP + FP + FN}$ | Fowlkes–Mallows index (FM)<br>$= \sqrt{PPV \times TPR}$        | Matthews correlation coefficient (MCC)<br>$= \frac{\sqrt{TPR \times TNR \times PPV \times NPV}}{\sqrt{FNR \times FPR \times FOR \times FDR}}$ | Threat score (TS), critical success index (CSI), Jaccard index<br>$= \frac{TP}{TP + FN + FP}$    |

# Other classification metrics based on **predicted classes**

Sources: [16][17][18][19][20][21][22][23] view · talk · edit

|                  |   | Predicted condition   |   |   |  |
|------------------|---|---|---|---|--|
|                  |   | Predicted positive (PP)   | Predicted negative (PN)   | Informedness, bookmaker informedness (BM)<br>$= \text{TPR} + \text{TNR} - 1$  | Prevalence threshold (PT)<br>$= \frac{\text{TP}}{\text{TPR} - \text{FPR}}$   |
| Actual condition | Total population<br>$= P + N$                                   |   |   |   |  |
|                  | Positive (P)<br><sup>[a]</sup>                                  | True positive (TP),<br>hit <sup>[b]</sup>   | False negative (FN),<br>miss,<br>underestimation                                    | True positive rate (TPR),<br>recall, sensitivity (SEN),<br>probability of detection, hit rate,<br>power<br>$= \frac{\text{TP}}{P} = 1 - \text{FNR}$   | False negative rate (FNR),<br>miss rate<br>type II error <sup>[c]</sup><br>$= \frac{\text{FN}}{P} = 1 - \text{TPR}$          |
|                  | Negative (N) <sup>[d]</sup>                                     | False positive (FP),<br>false alarm, overestimation   | True negative (TN),<br>correct rejection <sup>[e]</sup>                             | False positive rate (FPR),<br>probability of false alarm, fall-out<br>type I error <sup>[f]</sup><br>$= \frac{\text{FP}}{N} = 1 - \text{TNR}$   | True negative rate (TNR),<br>specificity (SPC), selectivity<br>$= \frac{\text{TN}}{N} = 1 - \text{FPR}$                      |
|                  | Prevalence<br>$= \frac{P}{P + N}$                               | Positive predictive value (PPV),<br>precision<br>$= \frac{\text{TP}}{\text{PP}} = 1 - \text{FDR}$   | False omission rate (FOR)<br>$= \frac{\text{FN}}{\text{PN}} = 1 - \text{NPV}$       | Positive likelihood ratio (LR+)<br>$= \frac{\text{TPR}}{\text{FPR}}$  | Negative likelihood ratio (LR-)<br>$= \frac{\text{FNR}}{\text{TNR}}$   |
|                  | Accuracy (ACC)<br>$= \frac{\text{TP} + \text{TN}}{P + N}$       | False discovery rate (FDR)<br>$= \frac{\text{FP}}{\text{PP}} = 1 - \text{PPV}$  | Negative predictive value (NPV)<br>$= \frac{\text{TN}}{\text{PN}} = 1 - \text{FOR}$ | Markedness (MK), deltaP ( $\Delta p$ )<br>$= \text{PPV} + \text{NPV} - 1$   | Diagnostic odds ratio (DOR)<br>$= \frac{\text{LR}+}{\text{LR}-}$   |
|                  | Balanced accuracy (BA)<br>$= \frac{\text{TPR} + \text{TNR}}{2}$ | F <sub>1</sub> score<br>$= \frac{2 \text{PPV} \times \text{TPR}}{\text{PPV} + \text{TPR}}$<br>$= \frac{2 \text{TP}}{2 \text{TP} + \text{FP} + \text{FN}}$ | Fowlkes–Mallows index (FM)<br>$= \sqrt{\text{PPV} \times \text{TPR}}$               | Matthews correlation coefficient (MCC)<br>$= \frac{\sqrt{\text{TPR} \times \text{TNR} \times \text{PPV} \times \text{NPV}}}{\sqrt{\text{FNR} \times \text{FPR} \times \text{FOR} \times \text{FDR}}}$ | Threat score (TS),<br>critical success index (CSI), Jaccard index<br>$= \frac{\text{TP}}{\text{TP} + \text{FN} + \text{FP}}$ |

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

↑  
Probability of being  
in class K

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = 1) \\ \vdots \\ Pr(\hat{y}_n = 1) \end{bmatrix}$$

For binary classification, this is typically rewritten as a column vector (not an  $n \times 2$  matrix)

# Output of classification methods

**True (observed)**  
class labels

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

**Predicted**  
class labels

$$\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix}$$

$$\hat{y}_i \in \{c_1, \dots, c_K\}$$

**Predicted**  
class probabilities

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = c_1) & \dots & Pr(\hat{y}_1 = c_K) \\ \vdots & \ddots & \vdots \\ Pr(\hat{y}_n = c_1) & \dots & Pr(\hat{y}_n = c_K) \end{bmatrix}$$

↑  
Probability of being  
in class 1

↑  
Probability of being  
in class K

$$\hat{\mathbf{p}} = \begin{bmatrix} Pr(\hat{y}_1 = 1) \\ \vdots \\ Pr(\hat{y}_n = 1) \end{bmatrix}$$

For binary classification, this is typically rewritten as a column vector (not an  $n \times 2$  matrix)

## Other classification metrics based on **predicted probabilities**

---

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels
  - Ex. For a binary classification problem, if  $P(y = 1) > 0.5$ , we typically predict class 1 and class 0 otherwise. However, if the harm of a false positive is very high, we might want to choose a different (higher) cutoff

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels
  - Ex. For a binary classification problem, if  $P(y = 1) > 0.5$ , we typically predict class 1 and class 0 otherwise. However, if the harm of a false positive is very high, we might want to choose a different (higher) cutoff
- + This leads to metrics like

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels
  - Ex. For a binary classification problem, if  $P(y = 1) > 0.5$ , we typically predict class 1 and class 0 otherwise. However, if the harm of a false positive is very high, we might want to choose a different (higher) cutoff
- + This leads to metrics like
  - **Area under the ROC (AUROC)**, or

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels
  - Ex. For a binary classification problem, if  $P(y = 1) > 0.5$ , we typically predict class 1 and class 0 otherwise. However, if the harm of a false positive is very high, we might want to choose a different (higher) cutoff
- + This leads to metrics like
  - **Area under the ROC (AUROC)**, or
  - **Area under the precision-recall curve (AUPRC)**

## Other classification metrics based on **predicted probabilities**

---

- + **Idea:** Different probability cutoffs/thresholds can be chosen to yield different predicted class labels
  - Ex. For a binary classification problem, if  $P(y = 1) > 0.5$ , we typically predict class 1 and class 0 otherwise. However, if the harm of a false positive is very high, we might want to choose a different (higher) cutoff
- + This leads to metrics like
  - **Area under the ROC (AUROC)**, or
  - **Area under the precision-recall curve (AUPRC)**

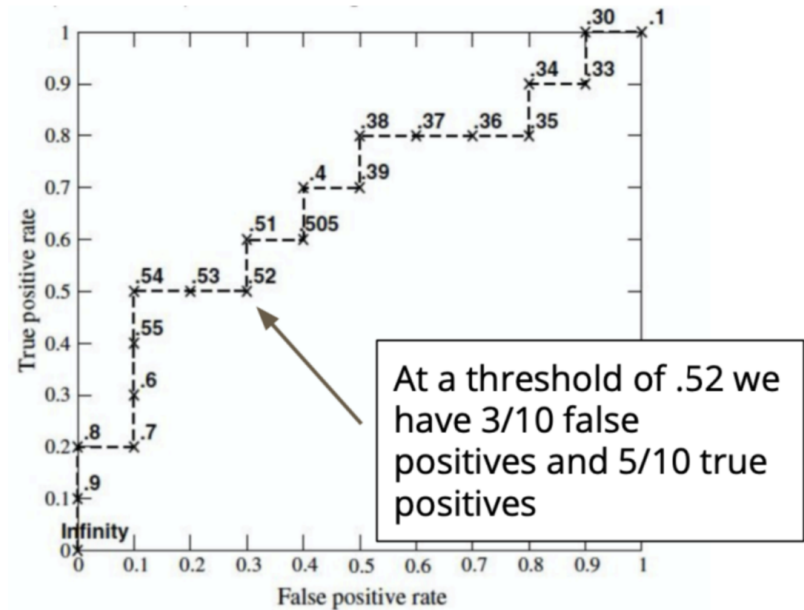
which summarize the classification accuracy across many different choices of thresholds

# Receiver operating characteristics (ROC) curve

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
- $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$

| Inst# | Class | Score | Inst# | Class | Score |
|-------|-------|-------|-------|-------|-------|
| 1     | p     | .9    | 11    | p     | .4    |
| 2     | p     | .8    | 12    | n     | .39   |
| 3     | n     | .7    | 13    | p     | .38   |
| 4     | p     | .6    | 14    | n     | .37   |
| 5     | p     | .55   | 15    | n     | .36   |
| 6     | p     | .54   | 16    | n     | .35   |
| 7     | n     | .53   | 17    | p     | .34   |
| 8     | n     | .52   | 18    | n     | .33   |
| 9     | p     | .51   | 19    | p     | .30   |
| 10    | n     | .505  | 20    | n     | .1    |

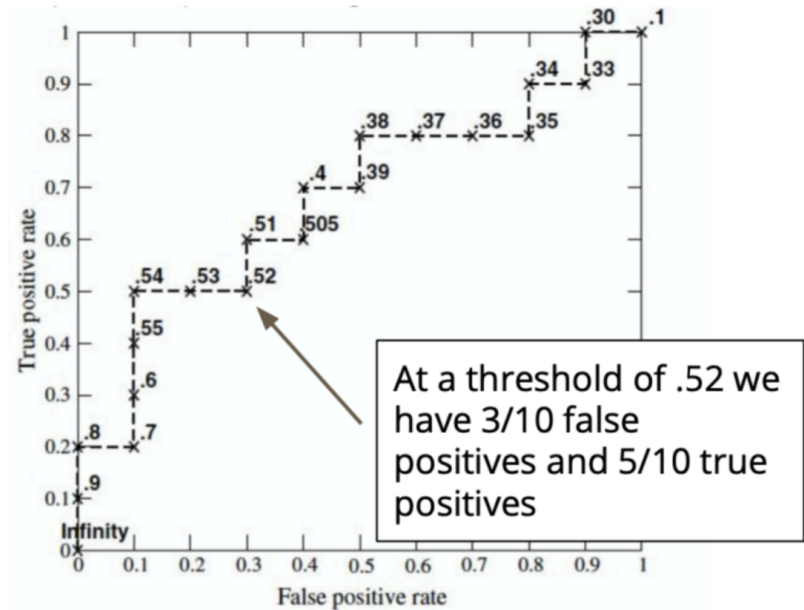


# Receiver operating characteristics (ROC) curve

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
- $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$

| Inst# | Class | Score | Inst# | Class | Score |
|-------|-------|-------|-------|-------|-------|
| 1     | p     | .9    | 11    | p     | .4    |
| 2     | p     | .8    | 12    | n     | .39   |
| 3     | n     | .7    | 13    | p     | .38   |
| 4     | p     | .6    | 14    | n     | .37   |
| 5     | p     | .55   | 15    | n     | .36   |
| 6     | p     | .54   | 16    | n     | .35   |
| 7     | n     | .53   | 17    | p     | .34   |
| 8     | n     | .52   | 18    | n     | .33   |
| 9     | p     | .51   | 19    | p     | .30   |
| 10    | n     | .505  | 20    | n     | .1    |



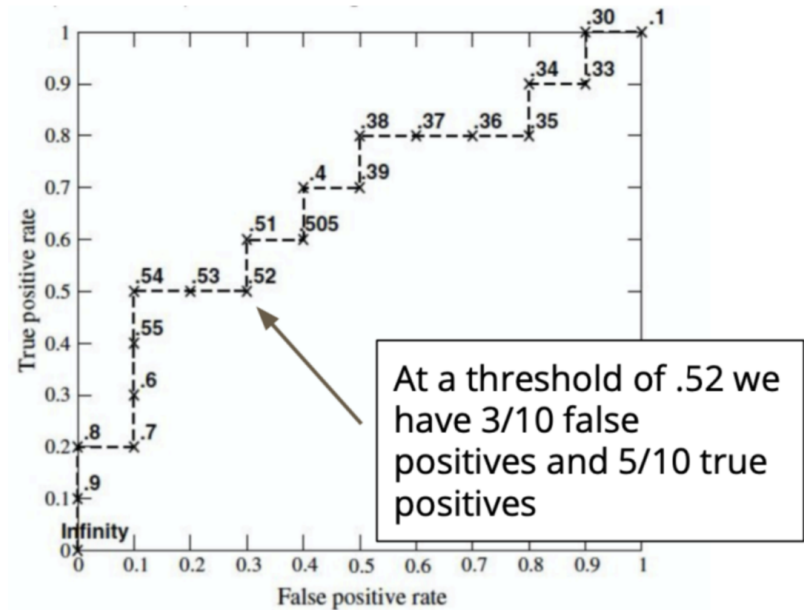
# Receiver operating characteristics (ROC) curve

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
- $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$

| Inst# | Class | Score | Inst# | Class | Score |
|-------|-------|-------|-------|-------|-------|
| 1     | p     | .9    | 11    | p     | .4    |
| 2     | p     | .8    | 12    | n     | .39   |
| 3     | n     | .7    | 13    | p     | .38   |
| 4     | p     | .6    | 14    | n     | .37   |
| 5     | p     | .55   | 15    | n     | .36   |
| 6     | p     | .54   | 16    | n     | .35   |
| 7     | n     | .53   | 17    | p     | .34   |
| 8     | n     | .52   | 18    | n     | .33   |
| 9     | p     | .51   | 19    | p     | .30   |
| 10    | n     | .505  | 20    | n     | .1    |

Classify  
positive





# Receiver operating characteristics (ROC) curve

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

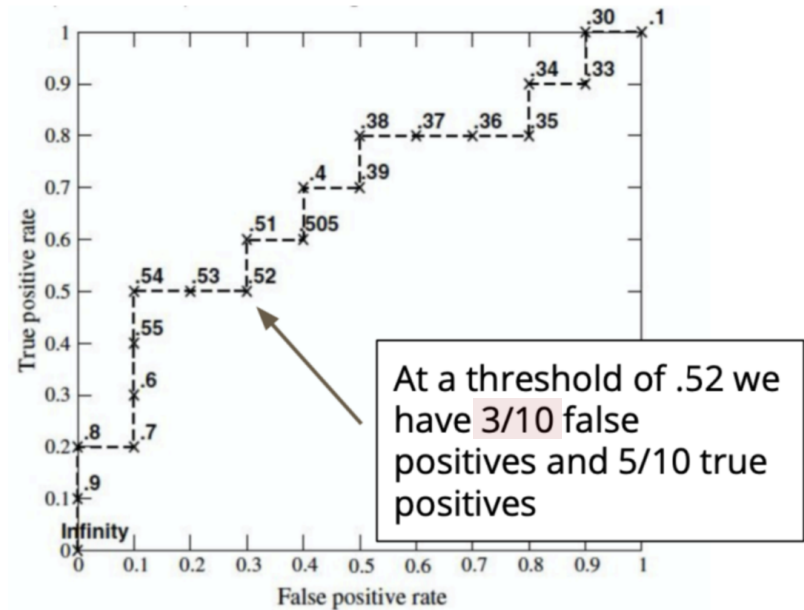
- $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
- $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$

| Inst# | Class | Score | Inst# | Class | Score |
|-------|-------|-------|-------|-------|-------|
| 1     | p     | .9    | 11    | p     | .4    |
| 2     | p     | .8    | 12    | n     | .39   |
| 3     | n     | .7    | 13    | p     | .38   |
| 4     | p     | .6    | 14    | n     | .37   |
| 5     | p     | .55   | 15    | n     | .36   |
| 6     | p     | .54   | 16    | n     | .35   |
| 7     | n     | .53   | 17    | p     | .34   |
| 8     | n     | .52   | 18    | n     | .33   |
| 9     | p     | .51   | 19    | p     | .30   |
| 10    | n     | .505  | 20    | n     | .1    |

Classify positive



Classify negative



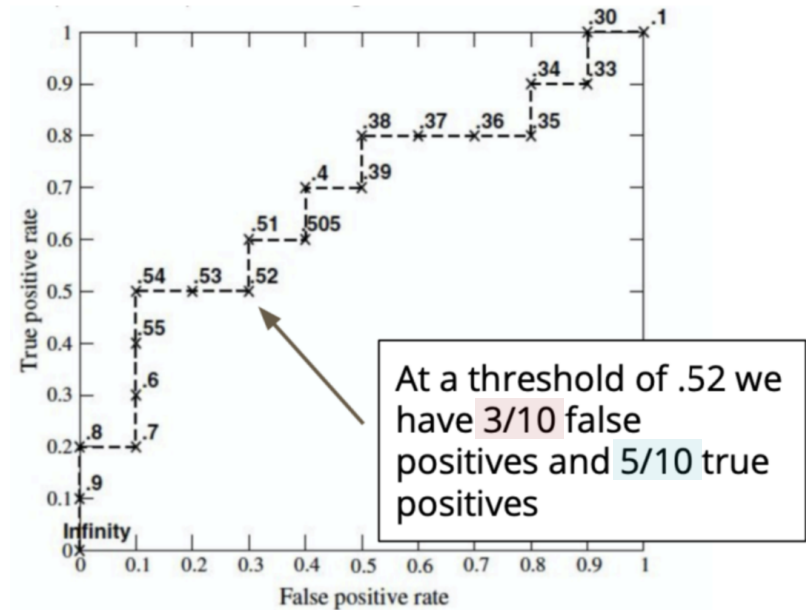
# Receiver operating characteristics (ROC) curve

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
- $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$

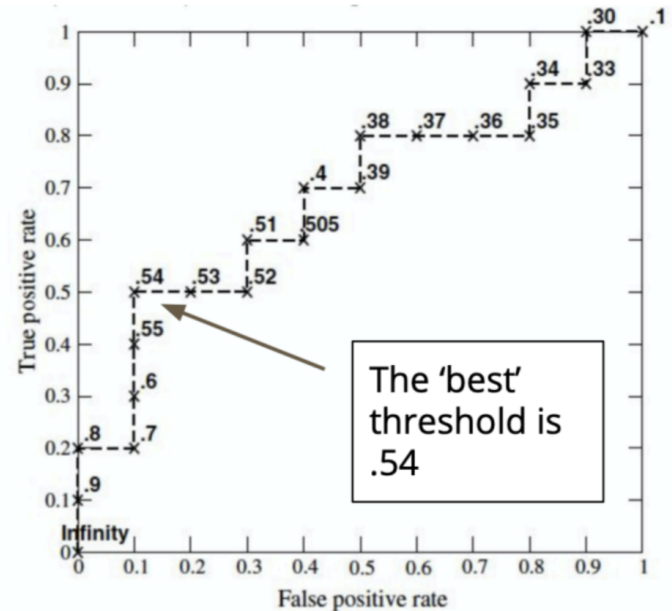
| Inst# | Class | Score | Inst# | Class | Score |
|-------|-------|-------|-------|-------|-------|
| 1     | p     | .9    | 11    | p     | .4    |
| 2     | p     | .8    | 12    | n     | .39   |
| 3     | n     | .7    | 13    | p     | .38   |
| 4     | p     | .6    | 14    | n     | .37   |
| 5     | p     | .55   | 15    | n     | .36   |
| 6     | p     | .54   | 16    | n     | .35   |
| 7     | n     | .53   | 17    | p     | .34   |
| 8     | n     | .52   | 18    | n     | .33   |
| 9     | p     | .51   | 19    | p     | .30   |
| 10    | n     | .505  | 20    | n     | .1    |

Classify  
positive ↑  
Classify  
negative ↓



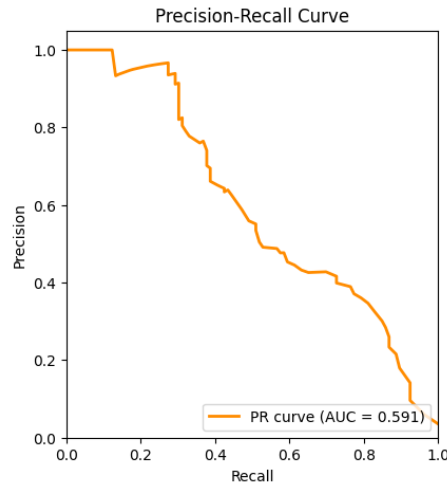
# Receiver operating characteristics (ROC) curve

- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)
  - $FPR = (\# \text{ false positives}) / (\# \text{ observed negative samples})$
  - $TPR = (\# \text{ true positives}) / (\# \text{ observed positive samples})$
- + The **area under the ROC curve (AUROC)** summarizes the results across all possible thresholds
- + The **“best” threshold** is often chosen to be the one closest to the top left corner of this plot (Youden's index)
  - Choose point which has the smallest distance to (0, 1)



# Precision-recall curve

- + **PR curve:** precision versus recall
  - Precision = (# true positives) / (# predicted positive samples)
  - Recall (aka TPR) = (# true positives) / (# observed positive samples)
- + The **area under the PR curve (AUPRC)** summarizes the results across all possible thresholds



# Comparing ROC and PR curves

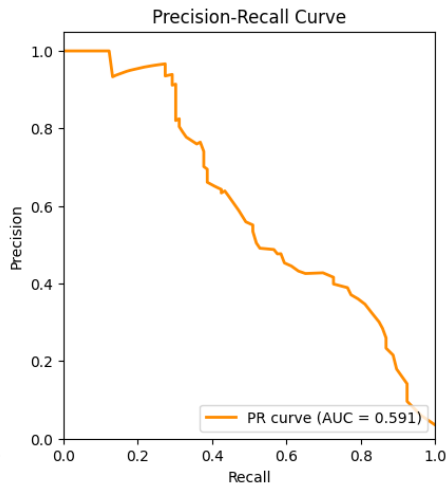
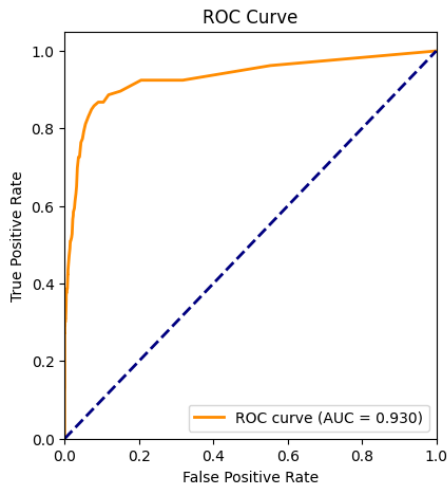
## + **PR curve:** precision versus recall

- Precision = (# true positives) / (# predicted positive samples)
- Recall (aka TPR) = (# true positives) / (# observed positive samples)

## + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- FPR = (# false positives) / (# observed negative samples)
- TPR (aka recall) = (# true positives) / (# observed positive samples)

AUROC  
baseline value = 0.5



AUPRC  
baseline value =  
proportion of  
samples in  
“positive” class

# Sidebar: Comparing "Baseline" ROC and PR curves

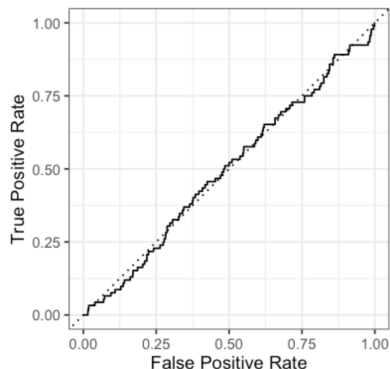
"Baseline" = predictions are completely independent of the true observed responses

**AUROC Baseline**  
= 0.5

**ROC**

$$\text{Recall/TPR} = \frac{\# \text{ True Positives}}{\# \text{ Observed Positives}}$$

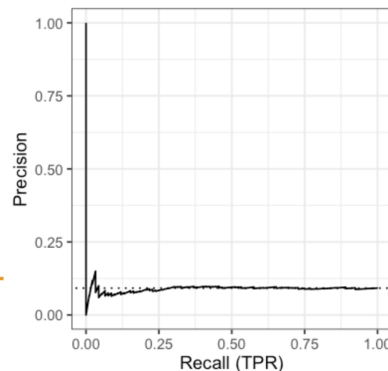
$$\text{FPR} = \frac{\# \text{ False Positives}}{\# \text{ Observed Negatives}}$$



**PR**

$$\text{Recall/TPR} = \frac{\# \text{ True Positives}}{\# \text{ Observed Positives}}$$

$$\text{Precision} = \frac{\# \text{ True Positives}}{\# \text{ Predicted Positives}}$$



**AUPRC Baseline = proportion of samples in positive class**

$$\text{TPR} = \frac{\pi \cdot PP}{\pi \cdot n}$$

$$= \frac{PP}{n}$$

$$\text{FPR} = \frac{(1 - \pi) \cdot PP}{(1 - \pi) \cdot n}$$

$$= \frac{PP}{n}$$

| Inst# | Class | Score |
|-------|-------|-------|
| 1     | N     | .9    |
| 2     | N     | .8    |
| 3     | P     | .7    |
| 4     | N     | .6    |
| 5     | N     | .55   |
| 6     | P     | .54   |
| 7     | N     | .53   |
| 8     | N     | .52   |
| 9     | N     | .51   |
| 10    | N     | .505  |

$PP$  { 3, 6 }  
 Classify positive ↑  
 Classify negative ↓

$$\text{TPR} = \frac{\pi \cdot PP}{\pi \cdot n}$$

$$= \frac{PP}{n}$$

$$\text{Precision} = \frac{\pi \cdot PP}{PP}$$

$$= \pi$$

where # samples =  $n$ , proportion of positive class =  $\pi$ , and # predicted positives =  $PP$

# Comparing ROC and PR curves

---

- + **PR curve:** precision versus recall
  - Precision = (# true positives) / (# predicted positive samples)
  - Recall (aka TPR) = (# true positives) / (# observed positive samples)
- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)
  - FPR = (# false positives) / (# observed negative samples)
  - TPR = (# true positives) / (# observed positive samples)

# Comparing ROC and PR curves

---

- + **PR curve:** precision versus recall
  - Precision = (# true positives) / (# predicted positive samples)
  - Recall (aka TPR) = (# true positives) / (# observed positive samples)
- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)
  - FPR = (# false positives) / (# observed negative samples)
  - TPR = (# true positives) / (# observed positive samples)
- + For highly imbalance problems where # negative (N) >> # positive (P) samples,

# Comparing ROC and PR curves

---

+ **PR curve:** precision versus recall

- Precision = (# true positives) / (# predicted positive samples)
- Recall (aka TPR) = (# true positives) / (# observed positive samples)

+ **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- FPR = (# false positives) / (# observed negative samples)
- TPR = (# true positives) / (# observed positive samples)

+ For highly imbalance problems where # negative (N) >> # positive (P) samples,

- FPR will always look pretty good no matter what FP is since FPR is dividing by N (a big number)

# Comparing ROC and PR curves

---

- + **PR curve:** precision versus recall

- Precision = (# true positives) / (# predicted positive samples)
- Recall (aka TPR) = (# true positives) / (# observed positive samples)

- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)

- FPR = (# false positives) / (# observed negative samples)
- TPR = (# true positives) / (# observed positive samples)

- + For highly imbalance problems where # negative (N) >> # positive (P) samples,
  - FPR will always look pretty good no matter what FP is since FPR is dividing by N (a big number)
  - Precision focuses instead on the predicted positive samples and avoids using N

# Comparing ROC and PR curves

---

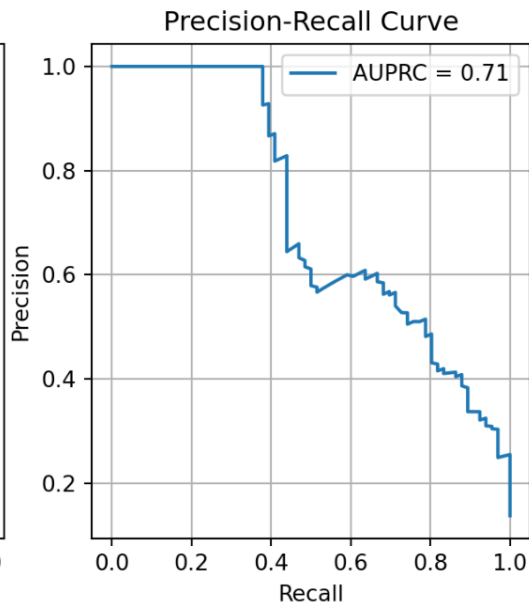
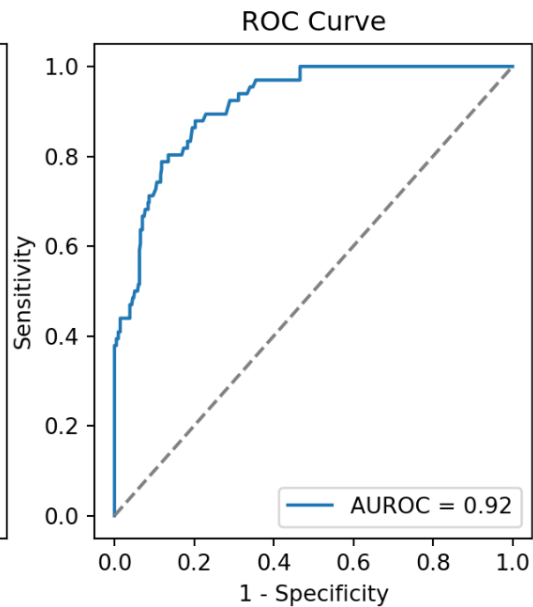
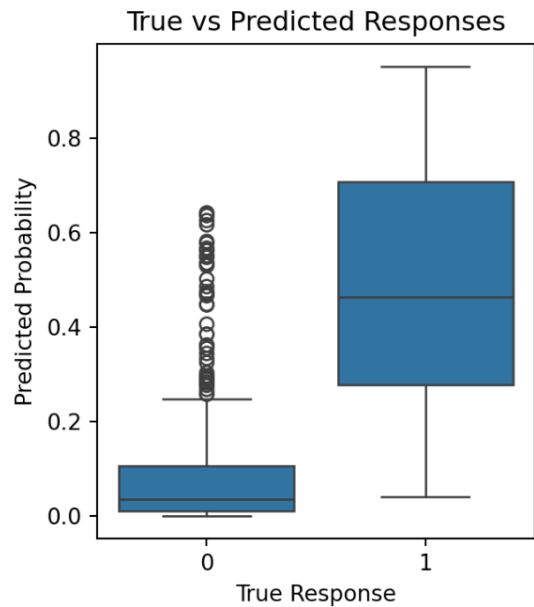
- + **PR curve:** precision versus recall
  - Precision = (# true positives) / (# predicted positive samples)
  - Recall (aka TPR) = (# true positives) / (# observed positive samples)
- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)
  - FPR = (# false positives) / (# observed negative samples)
  - TPR = (# true positives) / (# observed positive samples)
- + For highly imbalance problems where # negative (N) >> # positive (P) samples,
  - FPR will always look pretty good no matter what FP is since FPR is dividing by N (a big number)
  - Precision focuses instead on the predicted positive samples and avoids using N
- + As such, precision-recall curves are often preferred when there is high class imbalance

# Comparing ROC and PR curves

---

- + **PR curve:** precision versus recall
  - Precision = (# true positives) / (# predicted positive samples)
  - Recall (aka TPR) = (# true positives) / (# observed positive samples)
- + **ROC curve:** false positive rate (FPR) versus true positive rate (TPR)
  - FPR = (# false positives) / (# observed negative samples)
  - TPR = (# true positives) / (# observed positive samples)
- + For highly imbalance problems where # negative (N) >> # positive (P) samples,
  - FPR will always look pretty good no matter what FP is since FPR is dividing by N (a big number)
  - Precision focuses instead on the predicted positive samples and avoids using N
- + As such, precision-recall curves are often preferred when there is high class imbalance
- + **My general advice: always look at multiple metrics**

# Going beyond global summary metrics...



# Data Splitting

---

# Motivating Thought Experiment

---

**Ultimate question:** How well does my model perform on future, unseen data?

Consider the following workflow:

1. **Train/fit** prediction model on data  $\mathbf{X}$
2. **Evaluate** prediction model on same dataset  $\mathbf{X} \rightarrow$  demonstrates  $Z\%$  accuracy
3. **Generalizability:**  
“I expect my model to achieve  $Z\%$  accuracy when deployed for future data”

# Motivating Thought Experiment

---

**Ultimate question:** How well does my model perform on future, unseen data?

Consider the following workflow:

1. **Train/fit** prediction model on data  $\mathbf{X}$
2. **Evaluate** prediction model on same dataset  $\mathbf{X} \rightarrow$  demonstrates  $Z\%$  accuracy
3. **Generalizability:**  
“I expect my model to achieve  $Z\%$  accuracy when deployed for future data”

**Problem:** most likely over-fitting to the existing training data

# Overview of Data Splitting

---

**Data splitting is the key to assessing generalizability** (or how well our method performs on future unseen data)

# Overview of Data Splitting

---

**Data splitting is the key to assessing generalizability** (or how well our method performs on future unseen data)

The simplest case:

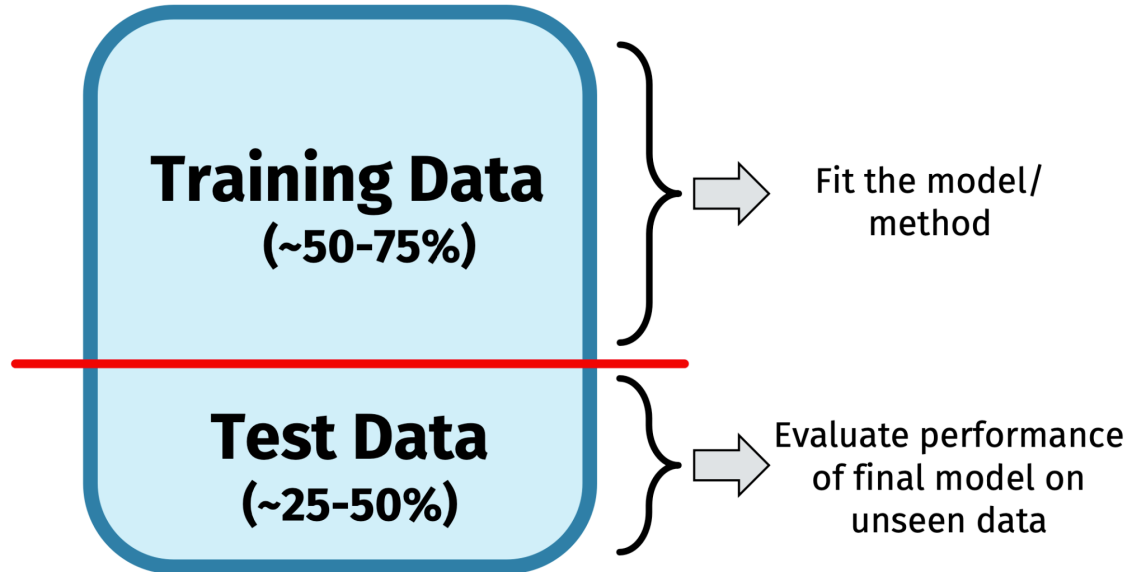
(ignoring choice of hyperparameters  
and possibility of multiple models)

# Overview of Data Splitting

**Data splitting is the key to assessing generalizability** (or how well our method performs on future unseen data)

The simplest case:

(ignoring choice of hyperparameters and possibility of multiple models)



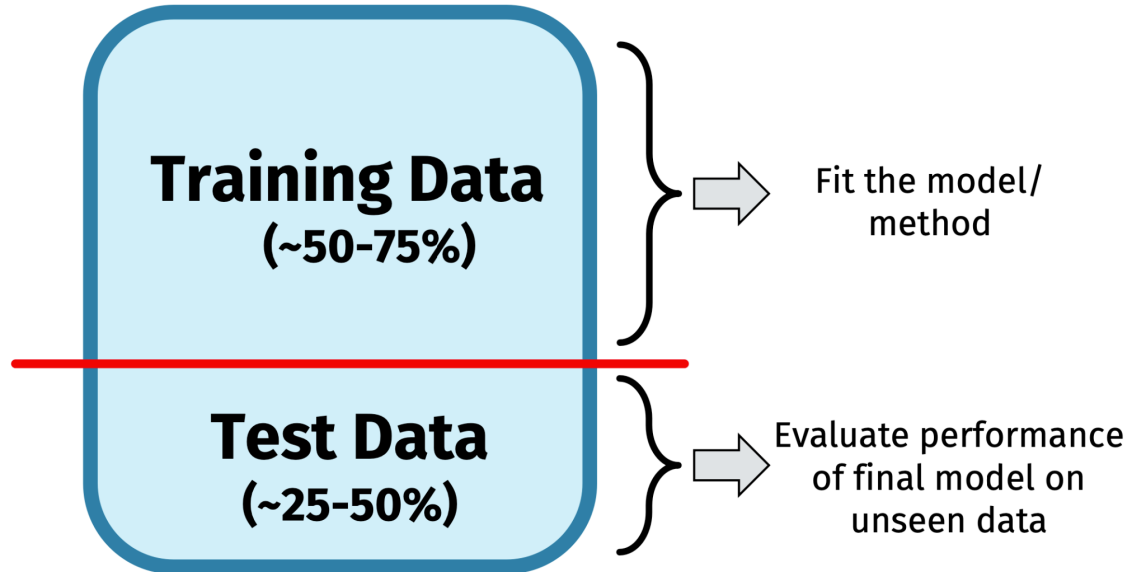
# Overview of Data Splitting

**Data splitting is the key to assessing generalizability** (or how well our method performs on future unseen data)

The simplest case:

(ignoring choice of hyperparameters and possibility of multiple models)

**Q:** How should we allocate samples to the training data versus the test data?



# How to choose data splitting scheme

---

**Q:** How should we allocate samples to the training data versus the test data?

- + Randomly? (default inclination)

Recall primary purpose of test set: evaluate generalizability to new data

↳ **Data splitting scheme should mimic the process of obtaining this future data**

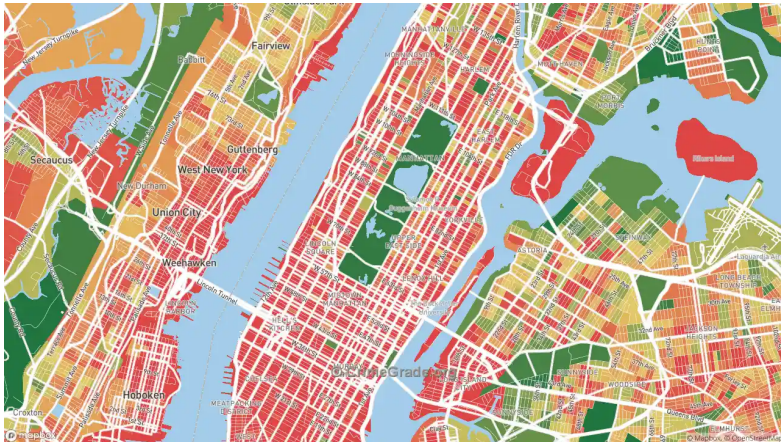
# How to choose data splitting scheme

---

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Spatial data:** Suppose we want to predict neighborhood crime rates and have data from NYC but want to apply/generalize model to LA



# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Spatial data:** Suppose we want to predict neighborhood crime rates and have data from NYC but want to apply/generalize model to LA



- Random splitting does not mimic the process of obtaining data from a *completely new region*
- + Need to **spatially block** data and put each block fully in train or fully in test

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Spatial data:** Suppose we want to predict neighborhood crime rates and have data from NYC but want to apply/generalize model to LA



- Random splitting does **not** mimic the process of obtaining data from a *completely new region*
- + Need to **spatially block** data and put each block fully in train or fully in test

# How to choose data splitting scheme

---

**Data splitting scheme should mimic the process of obtaining this future data**

**Q:** What would happen if we performed data splitting with random sampling and ignored the underlying process of obtaining future data?

**Our estimated test error would be **too optimistic****

- "Our results are too good to be true"

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Temporal (time) data:** Suppose we want to forecast K-days-ahead temperatures and we have historical data to use for training purposes

|        | temp. | month | ... |  |  |
|--------|-------|-------|-----|--|--|
| Time 1 |       |       |     |  |  |
|        |       |       |     |  |  |
| ⋮      |       |       |     |  |  |
|        |       |       |     |  |  |
| Time T |       |       |     |  |  |

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Temporal (time) data:** Suppose we want to forecast K-days-ahead temperatures and we have historical data to use for training purposes

|        | temp. | month | ... |  |  |
|--------|-------|-------|-----|--|--|
| Time 1 |       |       |     |  |  |
|        |       |       |     |  |  |
| ⋮      |       |       |     |  |  |
|        |       |       |     |  |  |
|        |       |       |     |  |  |
| Time T |       |       |     |  |  |

- Random splitting does not mimic the process of obtain/predicting on future temperature data
- + Use time 1 through T-K for training to predict temperature K-days ahead
  - If we need to do multiple data splits, use a "rolling window" approach

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Q:** What if we also had temperature data from different cities and wanted to generalize to existing and new cities and future time points?

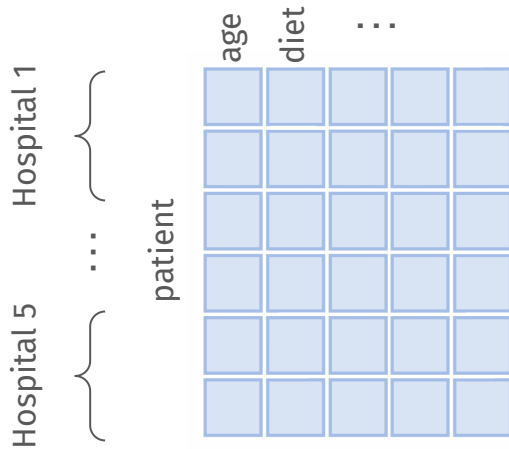
|                  | temp. | month | ... |  |  |
|------------------|-------|-------|-----|--|--|
| (city a, time 1) |       |       |     |  |  |
| ⋮                |       |       |     |  |  |
| (city a, time T) |       |       |     |  |  |
| (city b, time 1) |       |       |     |  |  |
| ⋮                |       |       |     |  |  |
| (city b, time T) |       |       |     |  |  |
| ⋮                |       |       |     |  |  |
|                  |       |       |     |  |  |

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Grouped/blocked data:** Suppose we want to develop a clinical model for risk of heart attacks and we have data from 5 hospitals

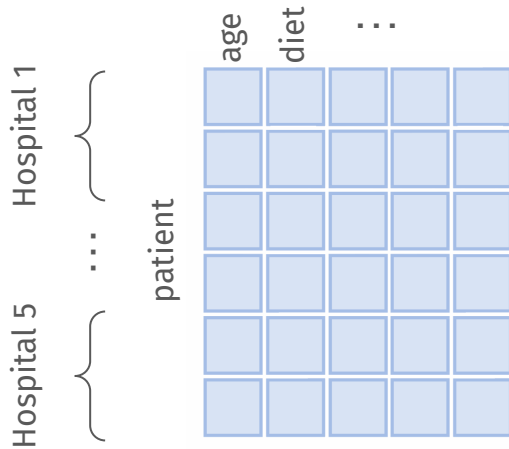


# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Grouped/blocked data:** Suppose we want to develop a clinical model for risk of heart attacks and we have data from 5 hospitals



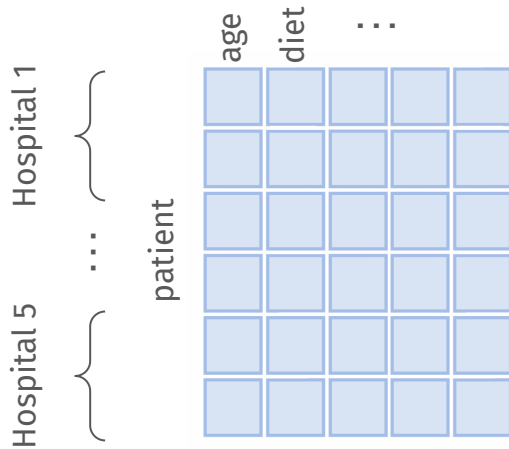
- + If we want to deploy our model just for those 5 hospitals, random splitting is ok

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Grouped/blocked data:** Suppose we want to develop a clinical model for risk of heart attacks and we have data from 5 hospitals



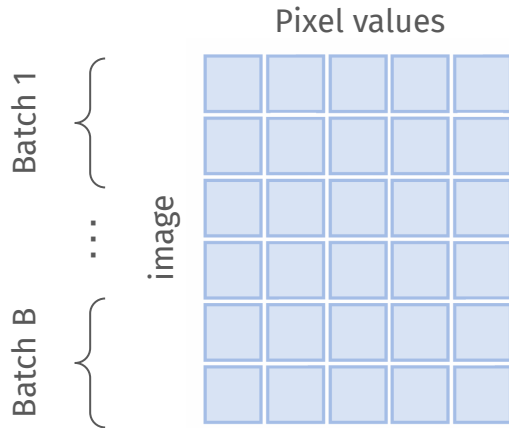
- + If we want to deploy our model just for those 5 hospitals, random splitting is ok
- + If we want to deploy our model to completely new hospitals, we should reserve at least one full hospital for the test set

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Grouped/blocked data:** Suppose we want to develop classify astronomical images from a particular telescope (or astronomer or experimental "batch")

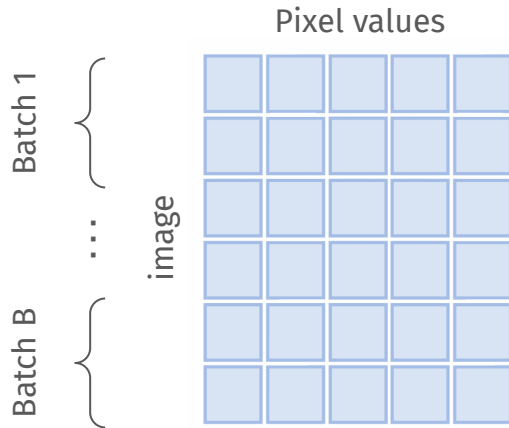


# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

- + **Grouped/blocked data:** Suppose we want to develop classify astronomical images from a particular telescope (or astronomer or experimental "batch")



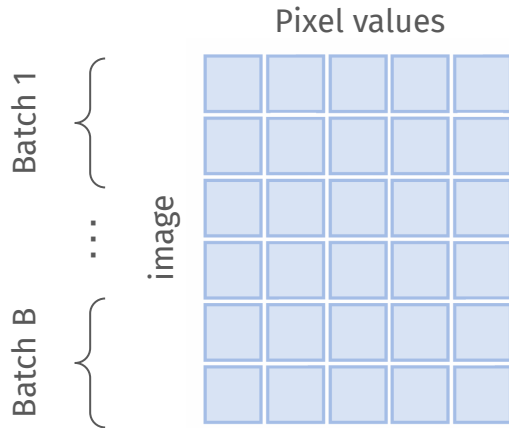
- + If we want to deploy our model just for that telescope/astronomer, random splitting is ok

# How to choose data splitting scheme

**Data splitting scheme should mimic the process of obtaining new (future) data**

*Examples:*

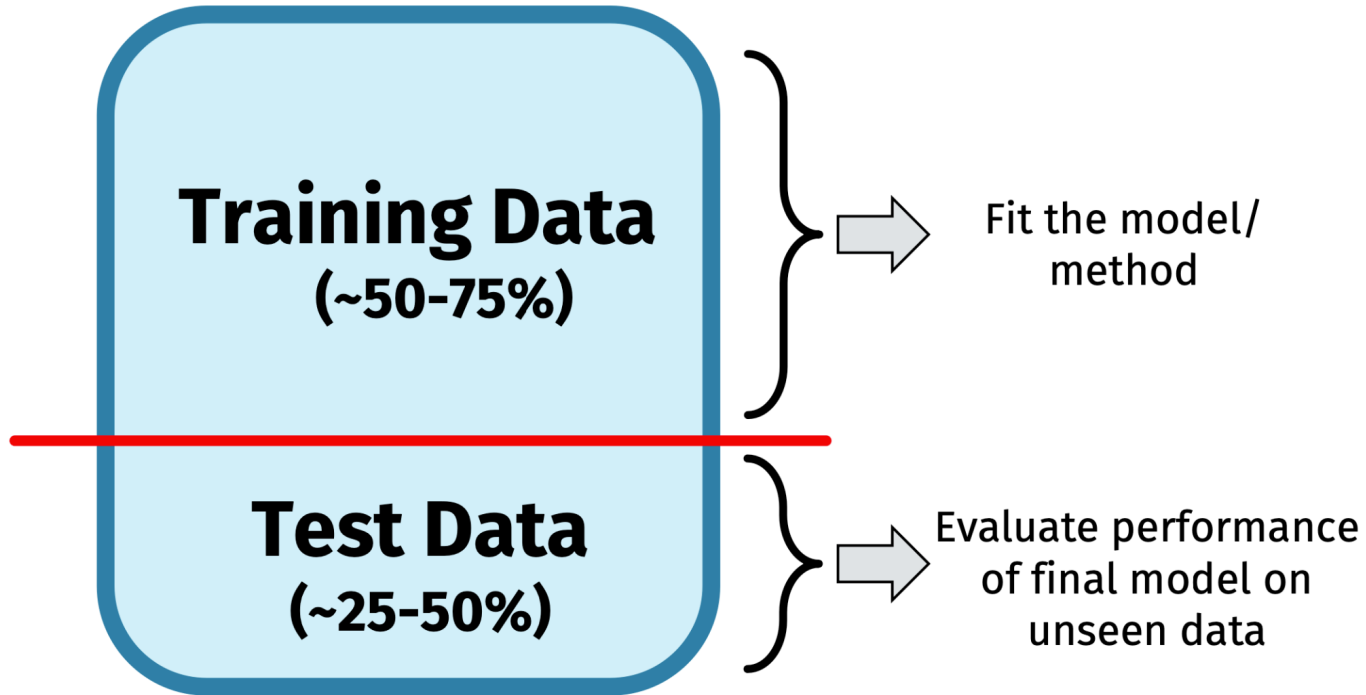
- + **Grouped/blocked data:** Suppose we want to develop classify astronomical images from a particular telescope (or astronomer or experimental "batch")



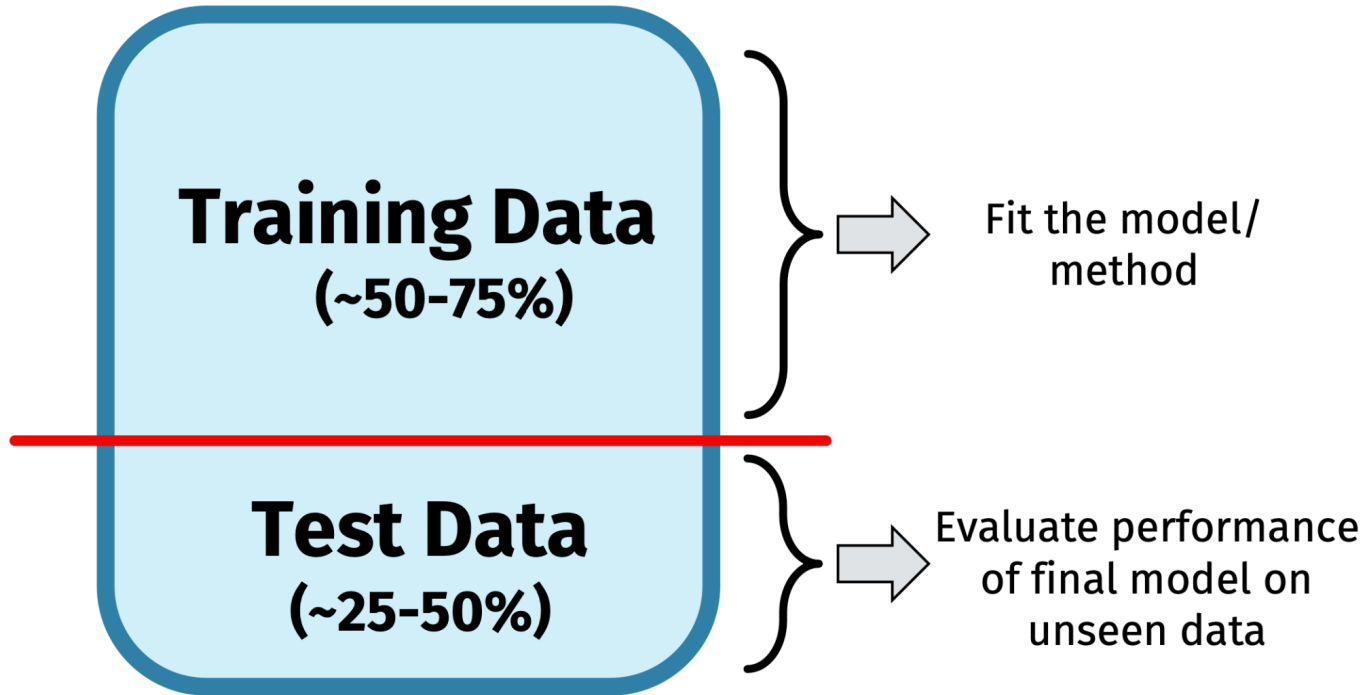
- + If we want to deploy our model just for that telescope/astronomer, random splitting is ok
- + If we want to deploy our model to a completely new telescope/astronomer, we should reserve at least one full "batch" for the test set

# Data splitting: the simplest case

---



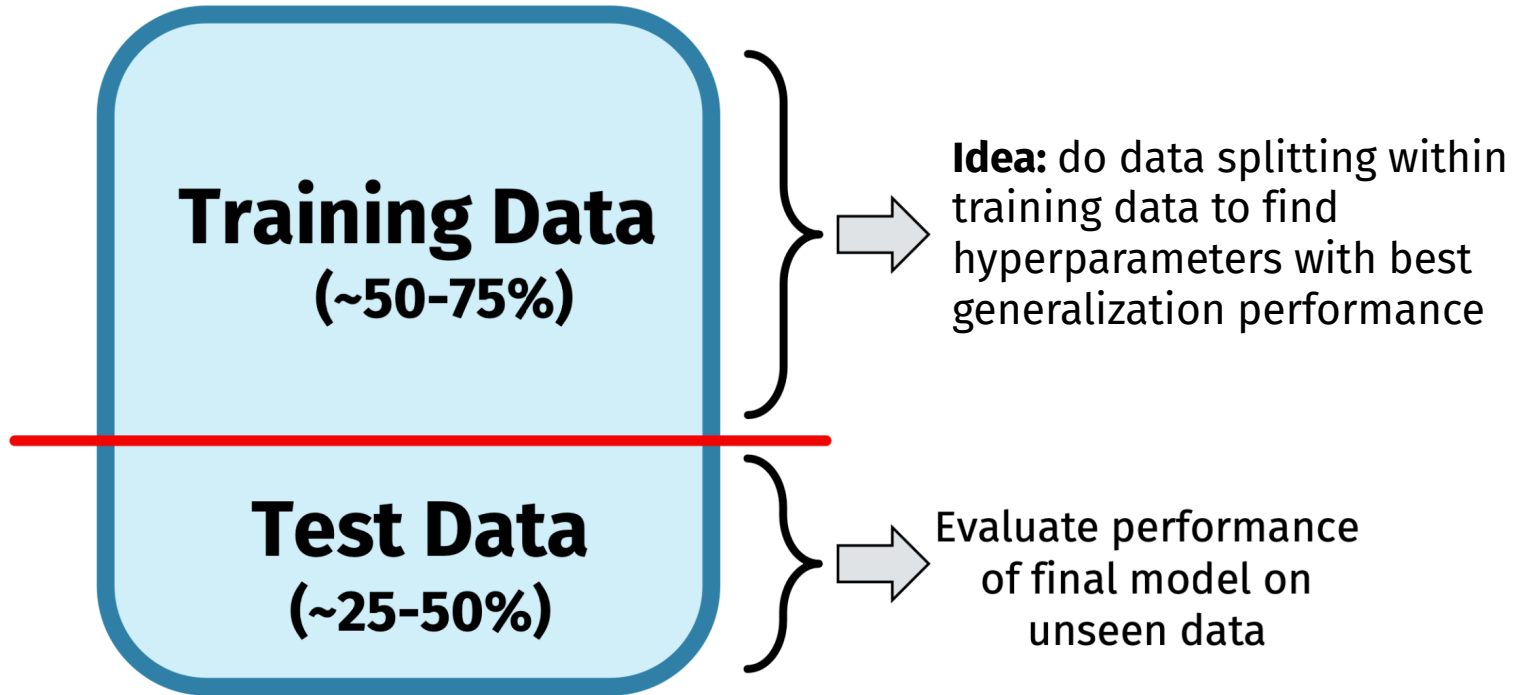
# Data splitting: the simplest case



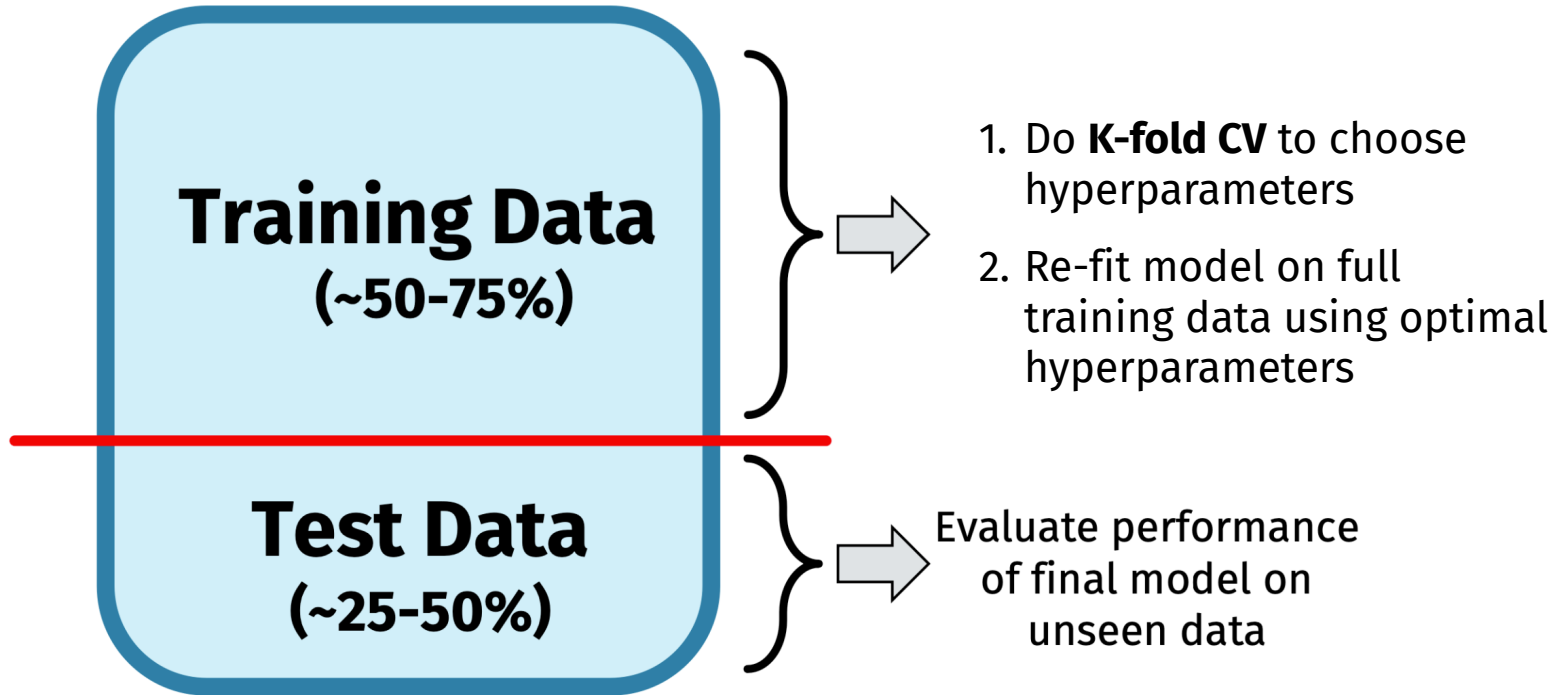
How can we modify this data splitting if we need to tune hyperparameters?

# Data splitting with hyperparameter tuning

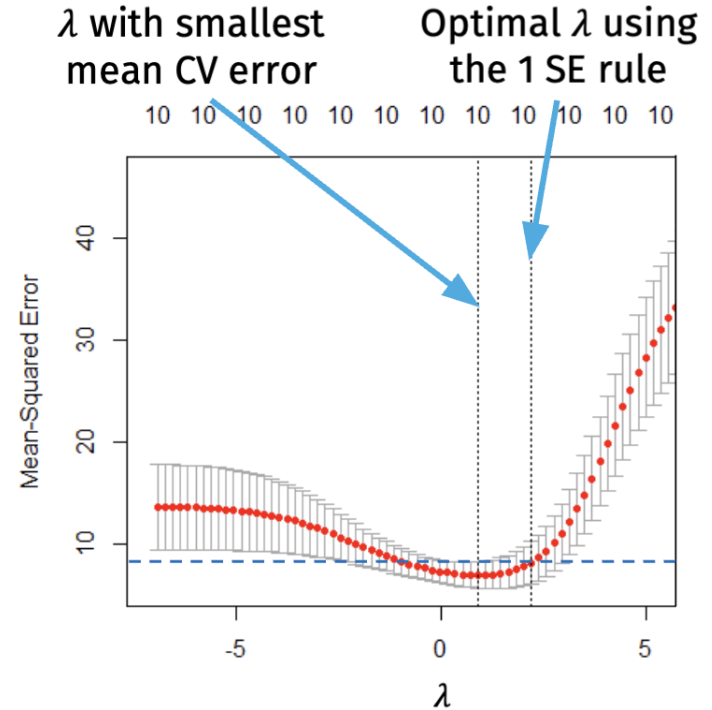
---



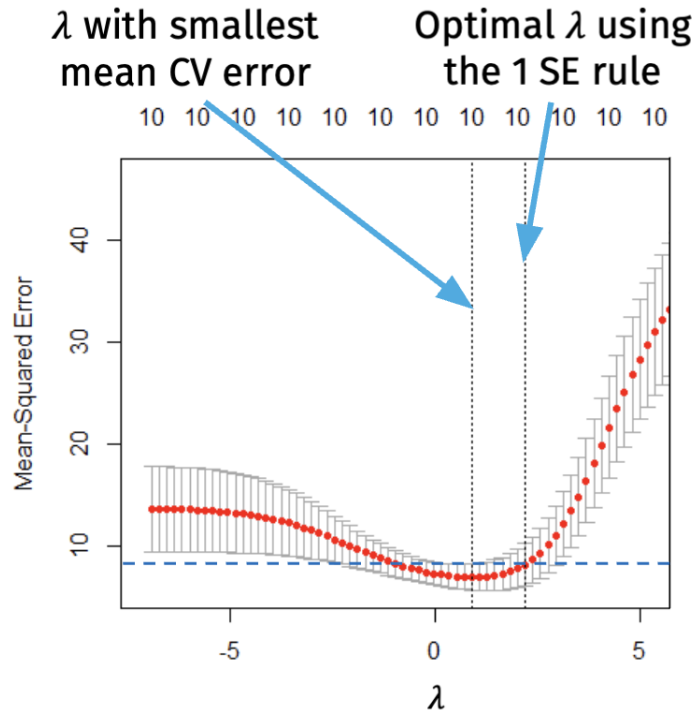
# Data splitting with hyperparameter tuning



# K-fold Cross-Validation (CV) for choosing hyperparameters



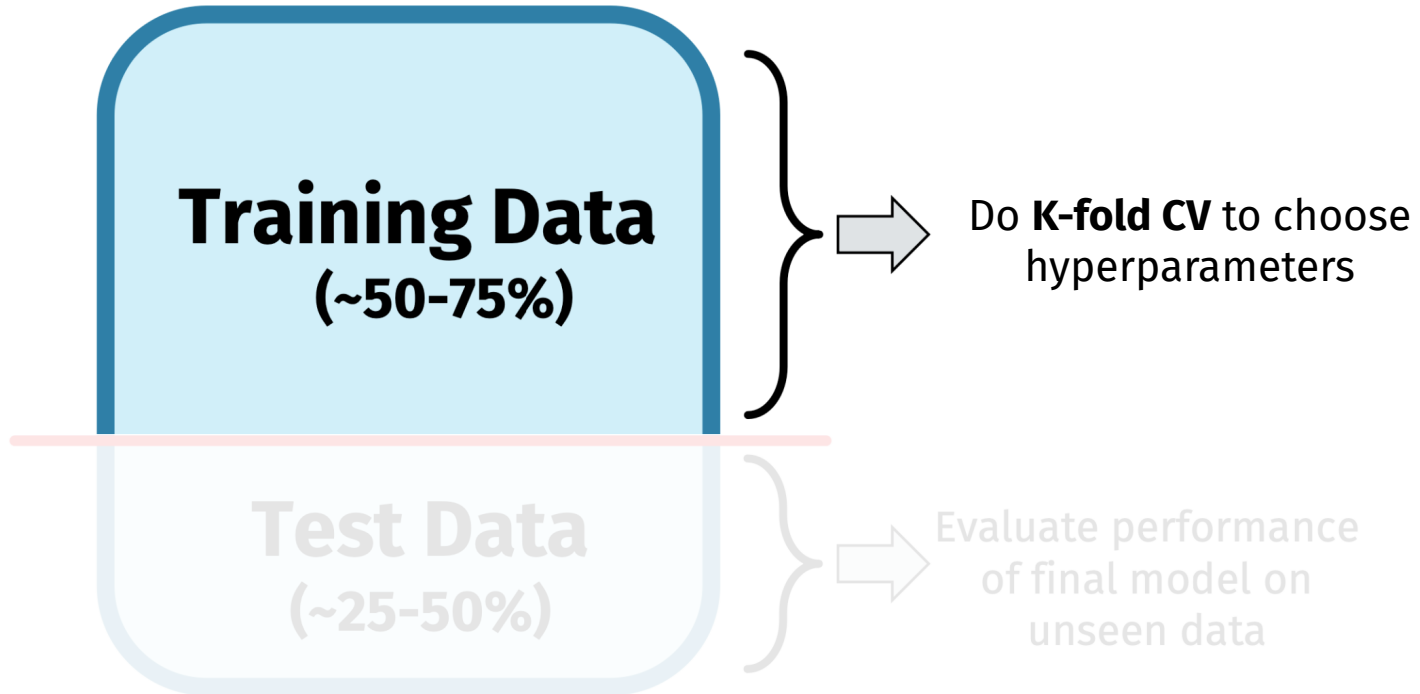
# K-fold Cross-Validation (CV) for choosing hyperparameters



**How do we split the data into folds?** Typically similar to how we split into train vs test

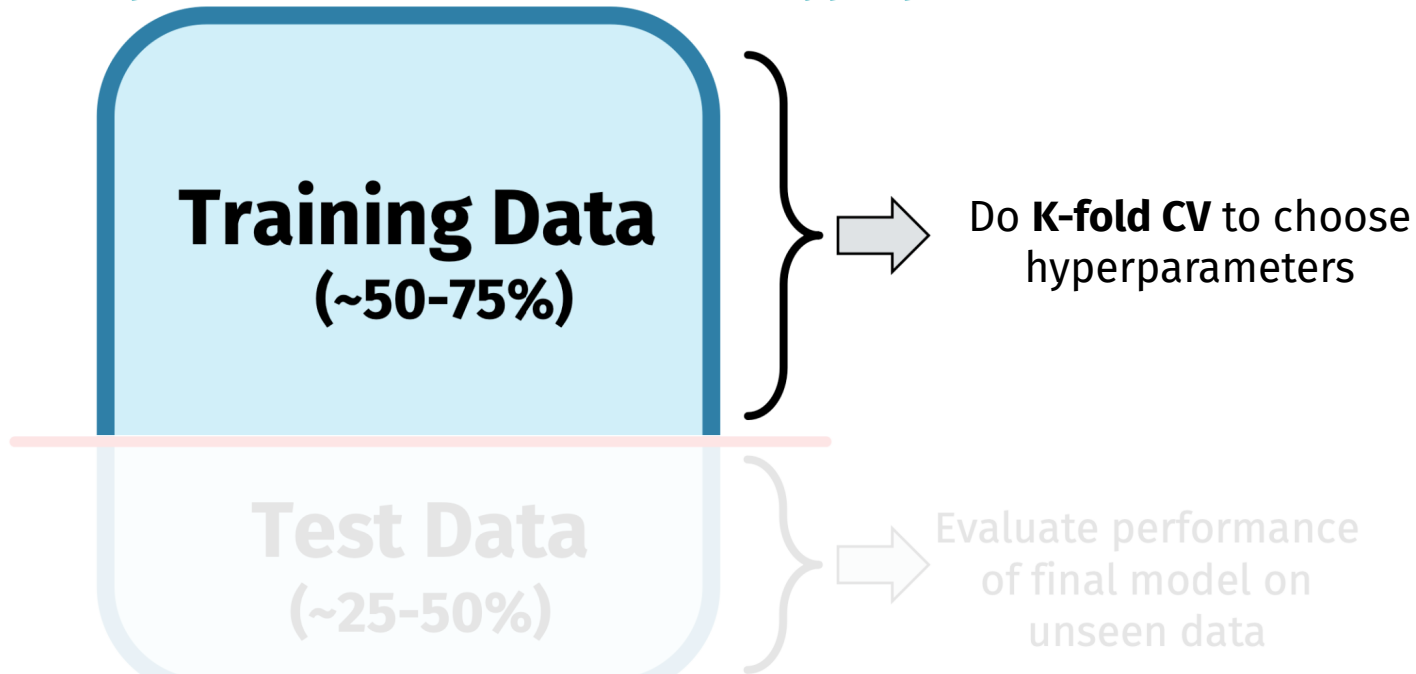
# Data splitting with hyperparameter tuning

What if we report the CV error for the best hyperparameter as our test error?



# Data splitting with hyperparameter tuning

What if we report the CV error for the best hyperparameter as our test error?

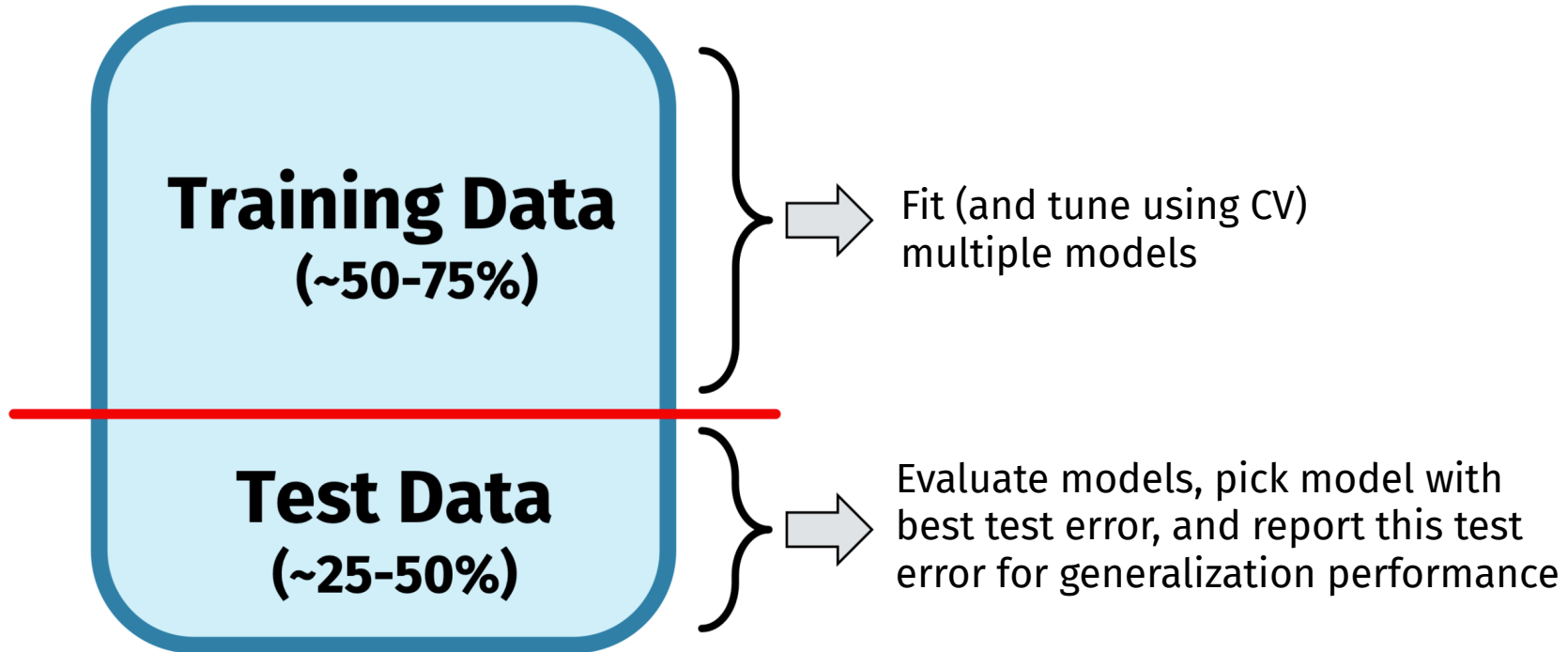


This would be an **overly-optimistic estimate** because we essentially took the minimum (or best) error across many attempts ("too good to be true")

# Data splitting with hyperparameter tuning + model selection

(multiple models to choose from)

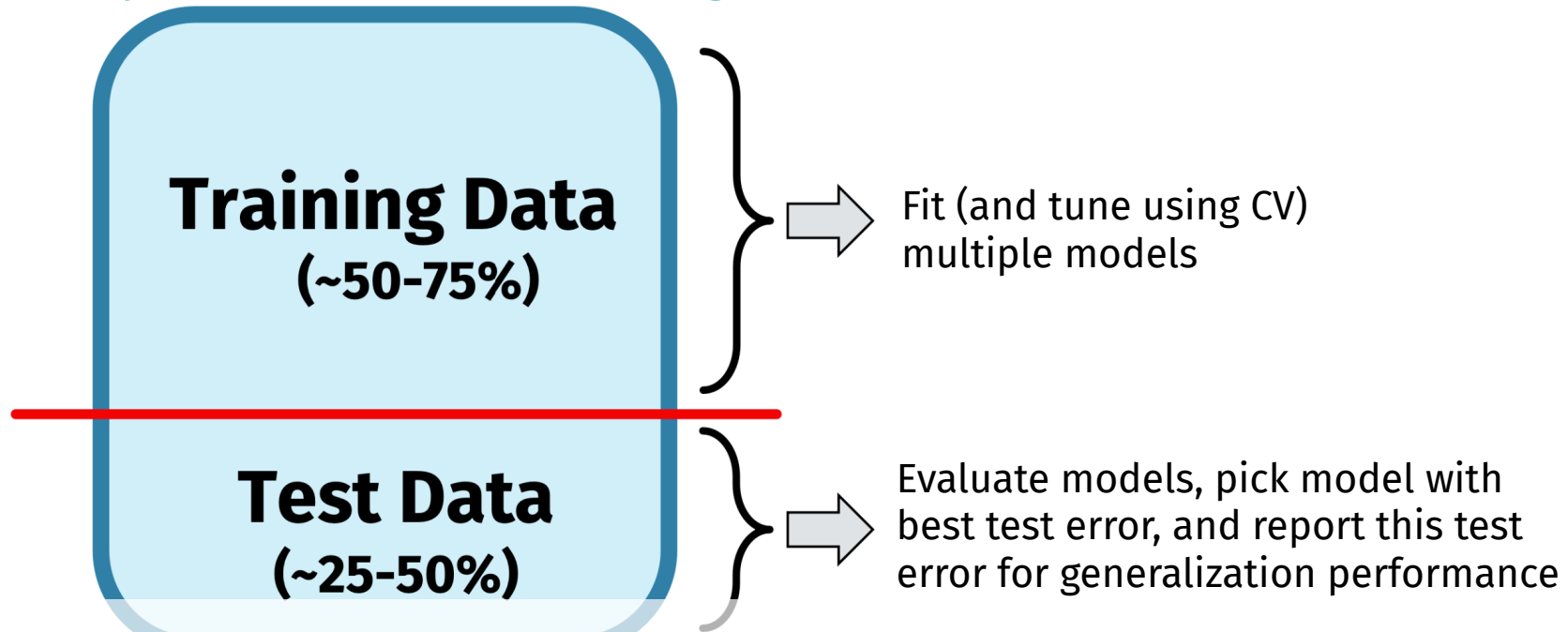
What if we report the test error after using test data to do model selection?



# Data splitting with hyperparameter tuning + model selection

(multiple models to choose from)

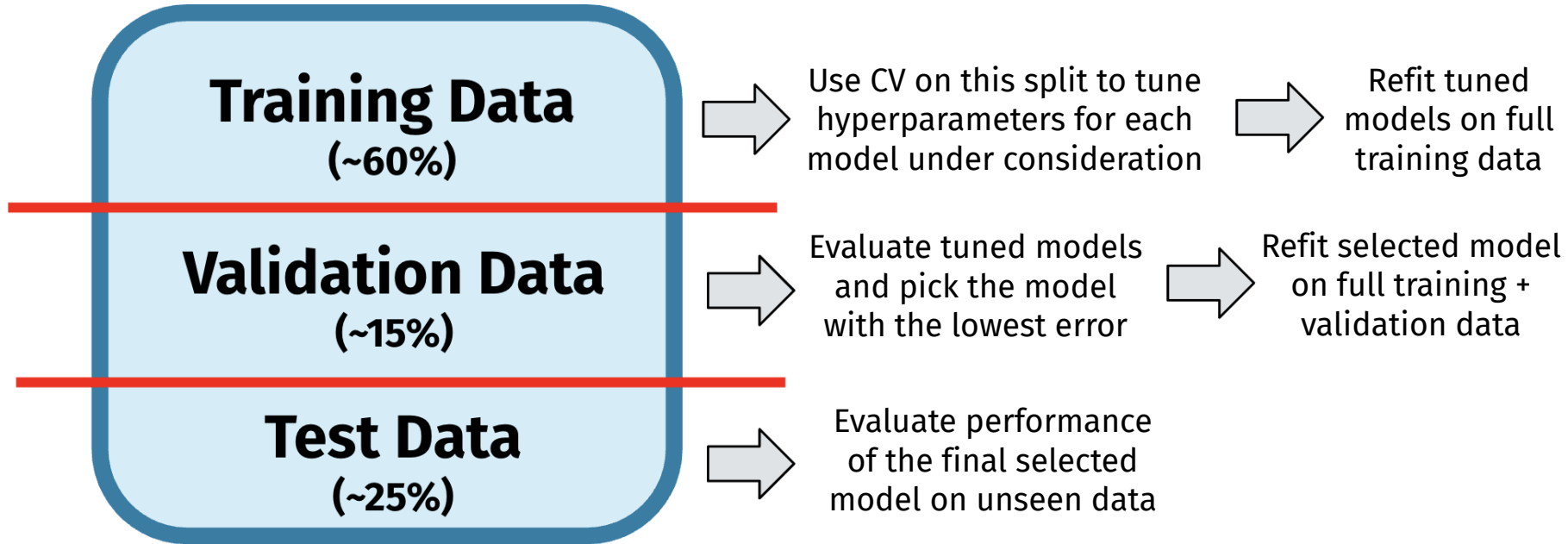
What if we report the test error after using test data to do model selection?



Again, this would be an **overly-optimistic estimate** because we essentially took the minimum (or best) error across many attempts ("too good to be true")

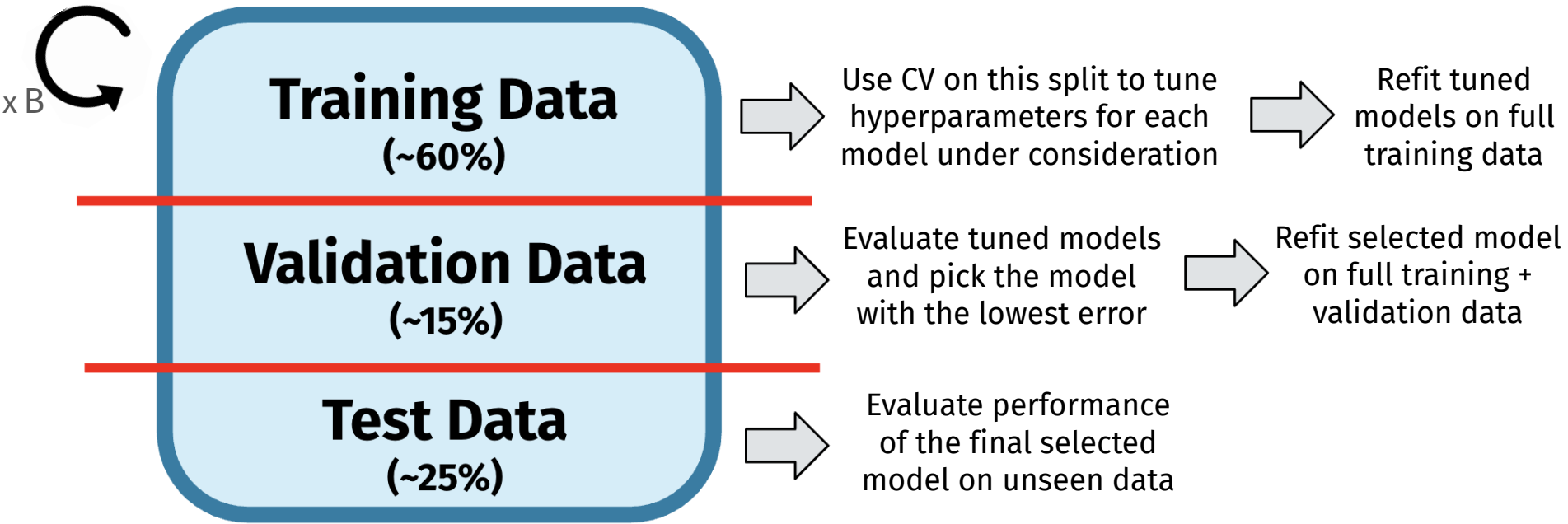
# Data splitting with hyperparameter tuning + model selection

(multiple models to choose from)



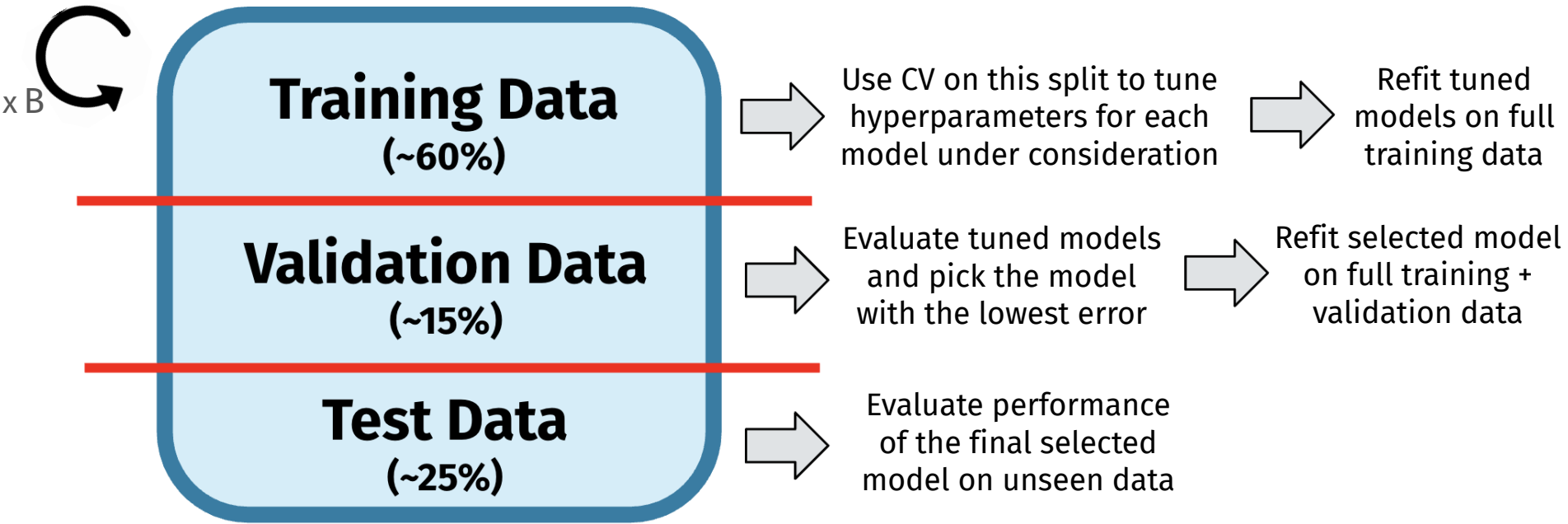
# Data splitting with hyperparameter tuning + model selection

(multiple models to choose from)



# Data splitting with hyperparameter tuning + model selection

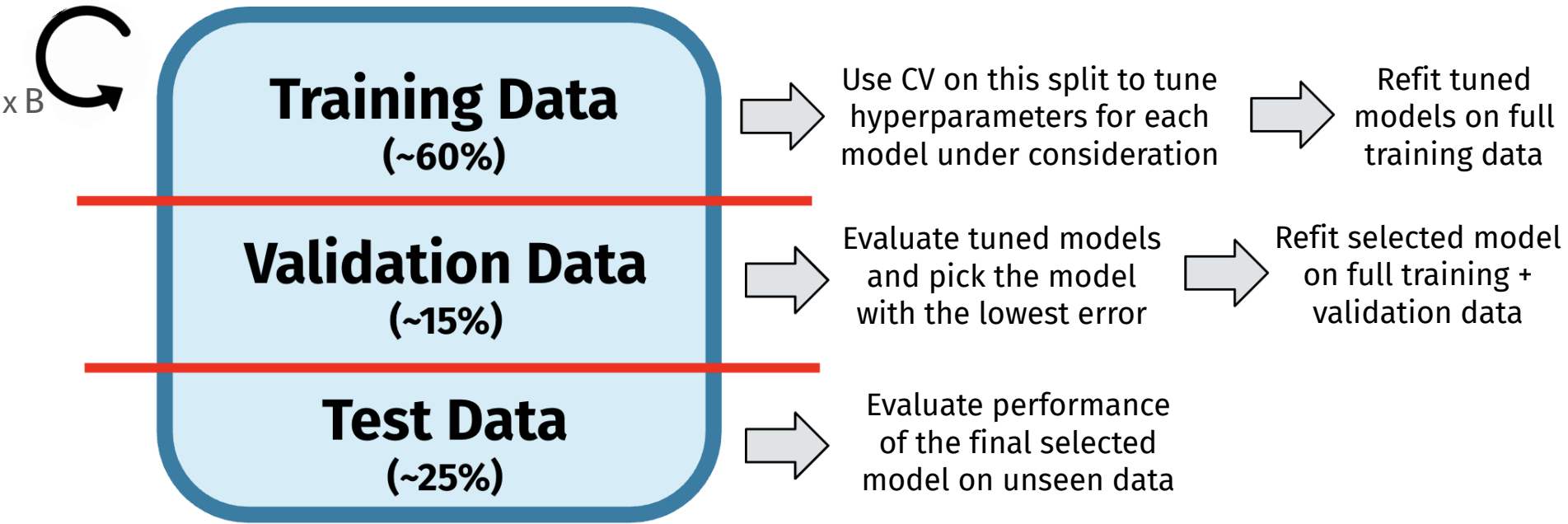
(multiple models to choose from)



- + Repeat this data splitting B times to get a variance estimate of the test error

# Data splitting with hyperparameter tuning + model selection

(multiple models to choose from)



- + Repeat this data splitting B times to get a variance estimate of the test error
- + This gives you an unbiased estimate of the prediction error for the **statistical learning pipeline/process**, NOT a specific model

# Key Takeaways of Data Splitting

---

## Data splitting is the key to assessing generalizability

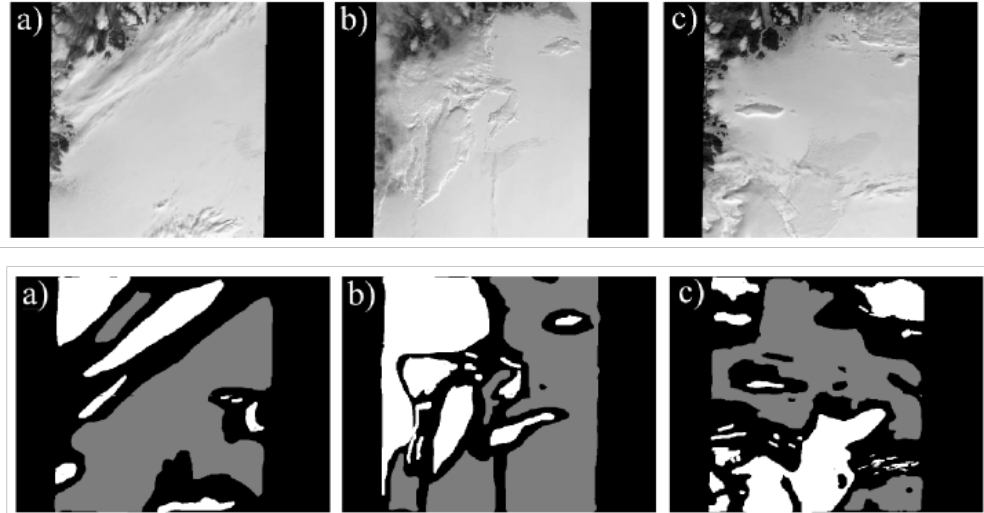
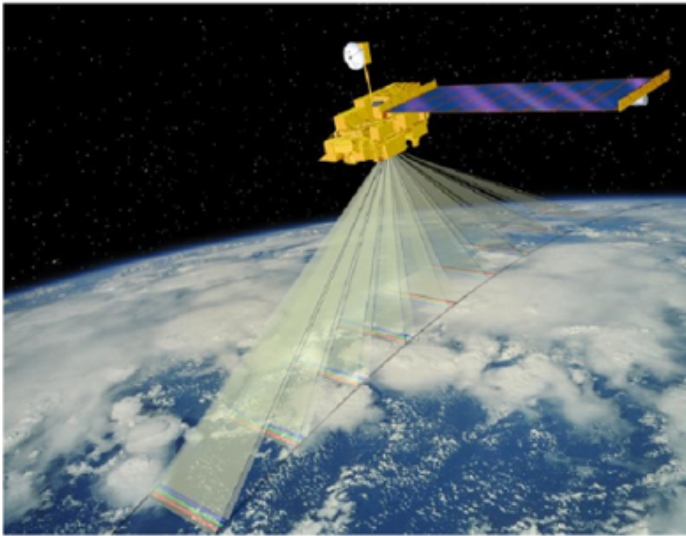
- + Need to decide how to (i) allocate samples into splits and (ii) splitting scheme
  - (i) Should mimic the process of obtaining new data in the future
  - (ii) Use training-validation-test split if selecting among multiple models;  
Use training-test if considering only one type of model
- + **Purpose of test set:** to obtain an unbiased estimate of the prediction error for your **statistical learning pipeline on future data**
  - Pipeline encompasses more than just the final prediction model. Includes the *process* of building your final prediction model
- + **IMPORTANT:** Do NOT touch the test set until the very end when you are ready to *evaluate* the generalization performance of your finalized pipeline

# Introduction to Lab 3

---

# Lab 3: Remote sensing for cloud detection (due March 6 6pm)

**Goal:** predict whether each pixel is a cloud or ice (from a glacier)



**Original study:** Shi et al. [Daytime Arctic Cloud Detection Based on Multi-Angle Satellite Data with Case Studies](#)