# Train / Validation / Test Split - 101

Bert Gollnick

# Train / Validation / Test Split - 101

| Training Data | Validation Data | Test Data |
|---|---|---|

- Data sample used to fit the model

# Train / Validation / Test Split - 101

| Training Data | Validation Data | Test Data |
|---|---|---|

- Data sample used to evaluate the model
- Used to fine-tune model hyperparameters
- Model occasionally „sees" the data, but does not learn from it
- Affects the model indirectly

# Train / Validation / Test Split - 101

| Training Data | Validation Data | Test Data |
|:---:|:---:|:---:|

- Data sample used to provide an unbiased evaluation of final model
- Provides gold-standard
- Model "sees" data only once
- used to evaluate competing models
- Same distribution as validation data
- Often only training and validation data is used, and no test data.

# Train / Validation / Test Split - 101

- Depends on two things: total number of samples, actual model
- Some models require more training data
- Validation data big enough to detect differences between models
- Models with few hyperparameters will be easy to validate → smaller validation dataset
- Models with many hyperparameters will be harder to validate → larger validation dataset

# Train / Validation / Test Split - 101

Interactive

Live Shiny-app!