

Machine Learning for Text Categorization

- **General setup:** Learn a classifier $f: X \rightarrow Y$
 - Input: X = all text objects; Output: Y = all categories
 - Learn a classifier function, $f: X \rightarrow Y$, such that $f(x)=y \in Y$ gives the correct category for $x \in X$ (“correct” is based on the training data)
- **All methods**
 - Rely on discriminative features of text objects to distinguish categories
 - Combine multiple features in a weighted manner
 - Adjust weights on features to minimize errors on the training data
- **Different methods** tend to vary in
 - Their way of measuring the errors on the training data (may optimize a different objective/loss/cost function)
 - Their way of combining features (e.g., linear vs. non-linear)