# Feature Selection with Mutual Information

Clayton W. Seitz

February 21, 2022

# Outline

References

# Feature Selection

What is it?

A special type of dimensionality reduction where we select a subset of features, in contrast with *feature extraction*

Why do we do it?

- ▶ Quality of the input data is just as important as the algorithm you choose
- ▶ The volume of a feature space grows exponentially in the number of dimensions $n$
- ▶ But we often have a small number of samples $p << n$

# Mutual Information

Mutual information comes from information theory and statistics.

$$I(X; Y) = D_{KL}(P(X, Y)||P(X)P(Y))$$
$$= H(X) - H(X|Y)$$

where $H$ denotes the entropy

- ▶ It quantifies the amount of information one variable carries about another
- ▶ Captures nonlinear correlation and is not limited to continuous variables
- ▶ $Y$ could be categorical e.g., cellular phenotypes

# An Example

# Using Mutual Information for Feature Selection

$$X^* = \underset{X}{\operatorname{argmax}}\ I(\boldsymbol{X}; Y)$$

For phenotyping, we might want to find the optimal $X$ which is most informative about the cell type $Y$

This is an optimization problem (NP-hard) on maximizing the *joint mutual information*
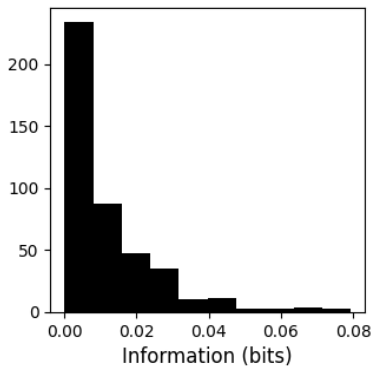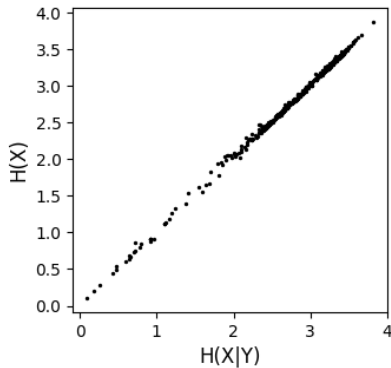
There are approximate solutions

# Maximum Relevancy Minimum Redundancy (MRMR)

By making some approximations, we can rewrite $I(\boldsymbol{X}; Y)$ as

$$I(\boldsymbol{X}; Y) \approx \sum_i \left( I(X_i; Y) - \alpha \sum_j I(X_i; X_j) \right)$$

where the parameter $\alpha$ determines how strongly we consider redundancy

# Results for T1D dataset

# Algorithm Details

The chain-rule for mutual information tells us that

$$I(\boldsymbol{X}; Y) = \sum_i I(X_i; Y | \boldsymbol{X}_{\setminus i}) \tag{1}$$

To simplify notation let $Z = \boldsymbol{X}_{\setminus i}$. The chain rule for info can also be used to show that

$$I(X; Y, Z) = I(X; Z) + I(X; Y | Z)$$

Solving for $I(X; Y | Z)$ says we can rewrite (1) as

$$I(\boldsymbol{X}; Y) = \sum_i I(X_i; Y | Z)$$
$$= \sum_i I(X_i; Y, Z) - I(X_i; Z)$$

# Algorithm Details

Applying the chain rule one more time gives

$$I(\boldsymbol{X}; Y) = \sum_i I(X_i; Y, Z) - I(X_i; Z)$$
$$= \sum_i I(X_i; Y) - I(X_i; Z) + I(X_i; Z|Y)$$

We maximize the sum by maximizing the each term $s_i$

$$s_i = I(X_i; Y) - I(X_i; Z) + I(X_i; Z|Y)$$
$$\approx I(X_i; Y) - \alpha \sum_j I(X_i; X_j) + \beta \sum_k I(X_i; X_k|Y)$$

Setting $\beta = 0$ gives the so-called maximum relevancy minimum redundancy (MRMR) features

## Algorithm Details

$$s_i \approx I(X_i; Y) - \alpha \sum_j I(X_i; X_j)$$

---
**Algorithm 1** Pseudocode for Greedy MRMR
---
1: *features* $= \{\}$
2: **for** $i = 1$ to $N$ **do**
3:    **if** $i = 1$ **then**
4:       add $x_i$ to features
5:    **else**
6:       **if** $s_i > s_{i-1}$ **then**
7:          add $x_i$ to features
8:       **end if**
9:    **end if**
10: **end for**
---

# References I