

# CS578 – INTERACTIVE AND TRANSPARENT MACHINE LEARNING

## TOPIC: FEATURE IMPORTANCE



**Mustafa Bilgic**



<http://www.cs.iit.edu/~mbilgic>



<https://twitter.com/bilgicm>

# TRANSPARENCY

- Which features are important?
- How important is each feature?
- How much does the importance of feature  $i$  decrease / increase after we add feature  $j$  to our domain?
- Is feature  $i$  positively or negatively correlated with the target feature?

# THREE MAIN APPROACHES

- Filter
  - Independent of the classification / regression model
  - E.g., mutual information
- Wrapper
  - Use the underlying model to select a subset of the features
  - E.g., could be done greedily bottom up or top down
- Joint optimization
  - Perform model training and feature selection jointly
  - E.g.,  $L_1$ -regularized logistic regression

# RELEVANT BUT DIFFERENT TOPIC

- Feature construction
- Difference
  - Feature selection *selects* a subset of features
  - Feature construction *constructs* a new feature space
- Examples of feature construction
  - Principal Component Analysis (PCA)
  - Factor Analysis (FA)
  - Linear Discriminant Analysis
  - See
    - <https://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>
    - [https://scikit-learn.org/stable/modules/classes.html#module-sklearn.discriminant\\_analysis](https://scikit-learn.org/stable/modules/classes.html#module-sklearn.discriminant_analysis)

# LET'S SEE A FEW EXAMPLES

- Now – filter
  - Information gain
  - Mutual information
  - $\text{Chi}^2$
- Later – subset and joint
  - Feature importances in decision trees
  - Feature weights in linear models
  - Regularization ( $L_1$  and  $L_2$ )
  - Influence analysis

# INFORMATION GAIN

$$\text{Entropy} = - \sum_j p_j \log_2 p_j$$

$$\text{InformationGain}(X_i) = \text{Entropy before } X_i - \text{Expected entropy after } X_i$$

# MUTUAL INFORMATION

$$MI(X, Y) = \sum_{x,y} \hat{P}(x, y) \log \frac{\hat{P}(x, y)}{\hat{P}(x)\hat{P}(y)}$$

# CHI<sup>2</sup>

$$\chi^2(X, Y) = \sum_{x, y} \frac{(M[x, y] - M \times \hat{P}(x) \times \hat{P}(y))^2}{M \times \hat{P}(x) \times \hat{P}(y)}$$



# REFERENCES

- JMLR article
  - <http://www.jmlr.org/papers/volume3/guyon03a/guyon03a.pdf>
- Scikit-learn
  - [https://scikit-learn.org/stable/modules/classes.html#module-sklearn.feature\\_selection](https://scikit-learn.org/stable/modules/classes.html#module-sklearn.feature_selection)
- Google Scholar
  - <https://scholar.google.com/scholar?q=feature+selection>
- Wikipedia
  - [https://en.wikipedia.org/wiki/Feature\\_selection](https://en.wikipedia.org/wiki/Feature_selection)

# YOUR TASK

- Read the references listed in the previous slide
- Get familiar with feature selection and ranking approaches on scikit-learn