# INFO251 – Applied Machine Learning
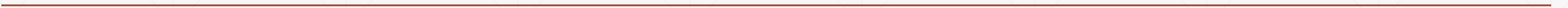
Lab 8 – Metrics & Fairness
**Satej Soman**
**based heavily on material by Suraj Nair, Joshua Blumenstock and Simón Ramirez Amaya**

# Announcements
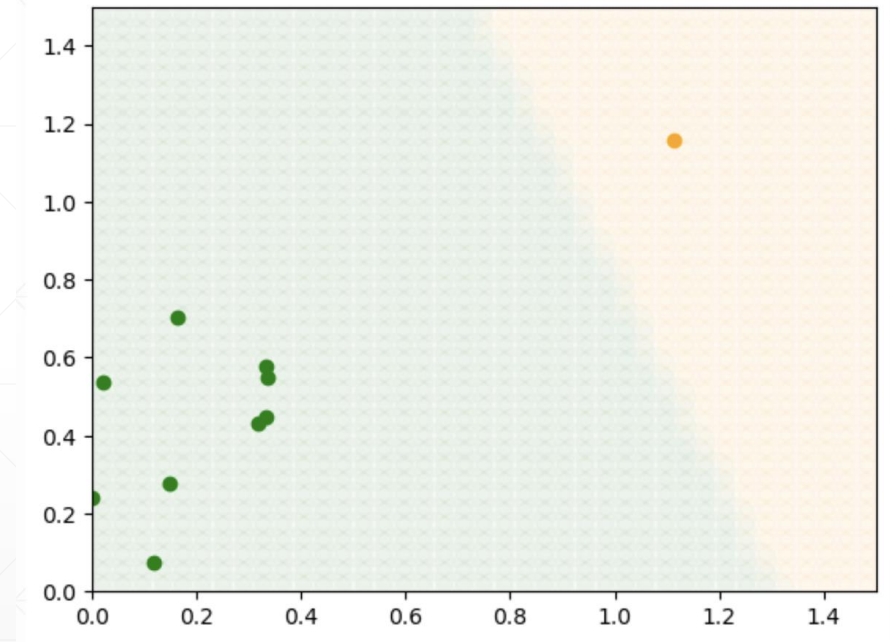
- PS5 due April 3

# Today's topics

- Review of prediction metrics

- Fairness:
  - Statistical non-discrimination
  - Fairness, datasets, benchmarks, and `folktables`
  - Useful components in `scikit-learn`

# Review: metrics for assessing prediction quality

- Assume that green is the "positive" class here

|        |        | Predicted | |
|--------|--------|-----------|--------|
|        |        | **Green** | **Orange** |
| **Actual** | **Green** | TP | FN |
|        | **Orange** | FP | TN |

- Accuracy = (TP + TN)/(TP + FP + FN + TN)

- TPR = TP/(TP + FN)

- FPR = FP/(FP + TN)

- Precision = TP/(TP + FP)

# Probabilistic interpretation

- Assume that green is the "positive" class here (Green = 1, Orange = 0)

- $\hat{y}$ is prediction, $y$ is the true value (both take on values 0 or 1)


- TPR = TP/(TP + FN) = $P(\hat{y} = 1 \mid y = 1)$

- FPR = FP/(FP + TN) = $P(\hat{y} = 1 \mid y = 0)$

- Precision = TP/(TP + FP) = $P(y = 1 \mid \hat{y} = 1)$

# More general metrics derived from *confusion matrix*

Source: https://en.wikipedia.org/wiki/Confusion_matrix

| | Predicted condition | | | |
|---|---|---|---|---|
| Total population = P + N | **Predicted positive** | **Predicted negative** | Informedness, bookmaker informedness (BM) $= TPR + TNR - 1$ | Prevalence threshold (PT) $= \frac{\sqrt{TPR \times FPR} - FPR}{TPR - FPR}$ |
| **Positive (P)** [a] | **True positive** (TP), hit[b] | **False negative** (FN), miss, underestimation | True positive rate (TPR), recall, sensitivity (SEN), probability of detection, hit rate, power $= \frac{TP}{P} = 1 - FNR$ | False negative rate (FNR), miss rate type II error [c] $= \frac{FN}{P} = 1 - TPR$ |
| **Negative (N)** [d] | **False positive** (FP), false alarm, overestimation | **True negative** (TN), correct rejection[e] | False positive rate (FPR), probability of false alarm, fall-out type I error [f] $= \frac{FP}{N} = 1 - TNR$ | True negative rate (TNR), specificity (SPC), selectivity $= \frac{TN}{N} = 1 - FPR$ |
| Prevalence $= \frac{P}{P + N}$ | Positive predictive value (PPV), precision $= \frac{TP}{TP + FP} = 1 - FDR$ | False omission rate (FOR) $= \frac{FN}{TN + FN} = 1 - NPV$ | Positive likelihood ratio (LR+) $= \frac{TPR}{FPR}$ | Negative likelihood ratio (LR−) $= \frac{FNR}{TNR}$ |
| Accuracy (ACC) $= \frac{TP + TN}{P + N}$ | False discovery rate (FDR) $= \frac{FP}{TP + FP} = 1 - PPV$ | Negative predictive value (NPV) $= \frac{TN}{TN + FN} = 1 - FOR$ | Markedness (MK), deltaP (Δp) $= PPV + NPV - 1$ | Diagnostic odds ratio (DOR) $= \frac{LR+}{LR-}$ |
| Balanced accuracy (BA) $= \frac{TPR + TNR}{2}$ | F₁ score $= \frac{2\, PPV \times TPR}{PPV + TPR} = \frac{2\, TP}{2\, TP + FP + FN}$ | Fowlkes–Mallows index (FM) $= \sqrt{PPV \times TPR}$ | Matthews correlation coefficient (MCC) $= \sqrt{TPR \times TNR \times PPV \times NPV} - \sqrt{FNR \times FPR \times FOR \times FDR}$ | Threat score (TS), critical success index (CSI), Jaccard index $= \frac{TP}{TP + FN + FP}$ |

(Actual condition, rows: Positive (P), Negative (N))
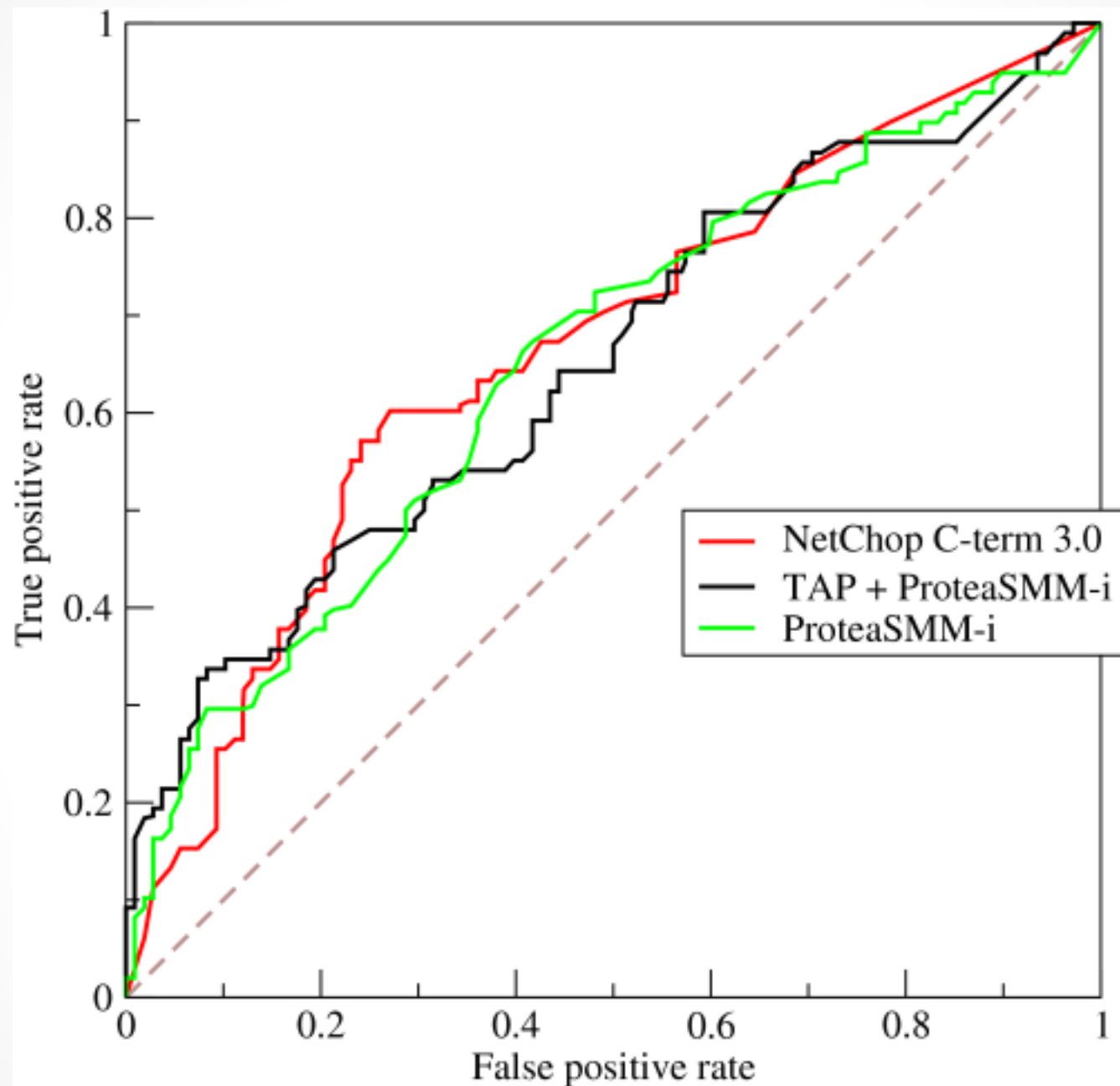
# Receiver Operator Characteristic Curve (ROC Curve)

- In classification problems, we usually get a predictive score that we threshold:
  - e.g. sigmoid predictor outputs a "score" in the range [0, 1]
  - a classifier takes this score, and compares it to a threshold
  - see: `model.predict_proba` in scikit-learn classifiers

- ROC Curve
  - Comes from WW2 – radar-based sensors scored by how well they detected enemy tanks/guns/ships
  - For every threshold, calculate the TPR and FPR and plot it out – see tradeoff between TPR and FPR
  - Other option for quota problems: Set "acceptance rate" to the rate of positive observations in the training set

# Hypothetical ROC Curves

# Example ROC Curves

Source:

https://en.wikipedia.org/wiki/Receiver_operating_characteristic

# Statistical definitions of fairness

- When we calculate TPR, FPR, etc – do we always calculate it for the entire dataset?

- The metrics we discussed (especially in their probabilistic formulation) are properties of the joint distribution of $y$ and $\hat{y}$

- What if some of the aspects of each data point also mattered? Consider a discrete variable $A$ (race, gender, location, etc) – it could be a feature, or it could be strongly correlated

- Now, the criteria are properties of the joint distribution of $y$, $\hat{y}$, and $A$

- **Error parity**: probability of error is the same for members and non-members
  - e.g. TPR parity: $P(\hat{y} = 1 \mid y = 1, A = 1) = P(\hat{y} = 1 \mid y = 1, A = 0)$

# Exercise on your own:

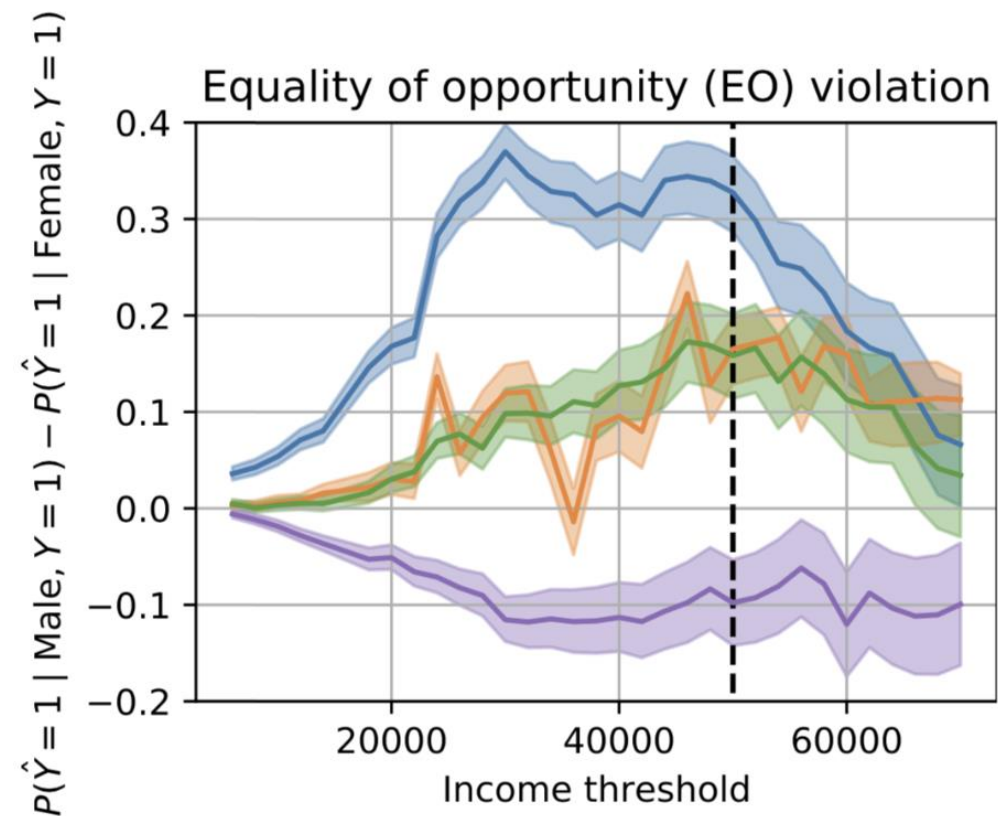- Does this classifier satisfy error parity for the TPR? For the FPR?
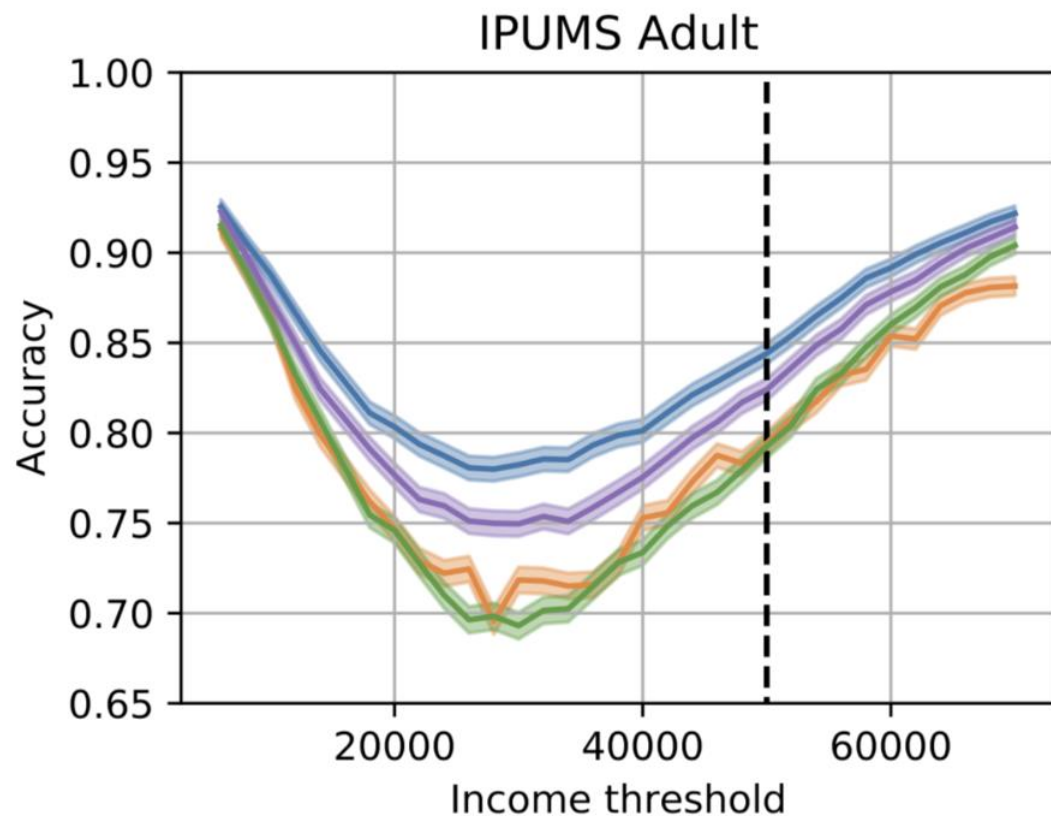
| A = 0 | $\hat{Y} = 0$ | $\hat{Y} = 1$ |
|-------|------|------|
| Y = 0 | 6/32 | 6/32 |
| Y = 1 | 1/32 | 3/32 |

| A = 1 | $\hat{Y} = 0$ | $\hat{Y} = 1$ |
|-------|------|------|
| Y = 0 | 4/32 | 4/32 |
| Y = 1 | 2/32 | 6/32 |

# Datasets and benchmarks

▪ Datasets are pretty key to ML practice!

▪ Performance on widely-used datasets (ImageNet, MNIST digits) becomes a benchmark to compare different methods' efficacy

▪ In ML-based fairness, the "UCI Adult" dataset was extremely popular – extract of 1994 US census data with a target binary variable based on whether income > $50k

▪ Hundreds if not thousands of ML fairness studies done using this dataset

▪ But… $50k was an arbitrary threshold! What if we used something else?

# Overfitting to UCI Adult

# folktables

- `folktables` is a Python package that offers a convenient and well documented interface to the ACS information!

- Using the resources in `folktables` you can generate data samples on the fly and get a better sense of whether a model is overfitting to a particular realization of the census data.

- For instance: you can train a model in CA data and used it to predict in AL. Or use data from two consecutive years to test how quickly performance degrades.

- Needless to say, you can include different covariates, set thresholds, etc...

https://github.com/socialfoundations/folktables

# i'm not directly concerned with fairness – is this relevant to general ML practice?

- yes.


- Understanding differential performance by subgroup can be
  - A. key to boosting overall performance
  - B. key to operationalizing predictions

# Useful components in sklearn

- Tree-based methods
  - sklearn.tree.DecisionTreeClassifier
  - sklearn.ensemble.RandomForestClassifier

- Pre-processing
  - sklearn.compose.ColumnTransformer
  - sklearn.preprocessing.OneHotEncoder
  - sklearn.preprocessing.StandardScaler

- Validation
  - sklearn.pipeline.Pipeline
  - sklearn.model selection.GridSearchCV
  - sklearn.model selection.cross validate