Ensembling $\rightarrow$ combine different models into 1

**Bagging**
$$y, \vec{X}$$

$$y_1, X_1 \quad y_2, X_2 \quad y_3, X_3$$

$$f_1 \quad f_2 \quad f_3$$

$$f(X) = g(f_1, f_2, f_3)$$

**Boosting**
$$y, \vec{X}$$

$$f_1(X) \underset{\sim}{\sim} y$$

$$e_1 = y - f_1$$

$$f_2(X) \underset{\sim}{\sim} e_1$$

$$e_2 = e_1 - f_1$$

$$f_3(X) \underset{\sim}{\sim} e_2$$

$$f = g(f_1, f_2, f_3)$$

**Stacking**
$$y, X$$

$$f_1(X) \quad f_2(X) \quad f_3(X)$$

$$f = g(f_1, f_2, f_3)$$

# Which model to use??



scikit-learn algorithm cheat-sheet

# General approaches to interpretable ML

**Without Machine Learning** — VERY SPECIFIC INSTRUCTIONS

**With Machine Learning** — DATA

Why did you predict 42 for this data point? *awkward silence*

Reliability $\Rightarrow$ small $\Delta x_i$ shouldn't produce large $\Delta \hat{y}_i$

Causality $\Rightarrow$ as we change $x_i$, can we anticipate $\Delta \hat{y}_i$

Trust $\Rightarrow$ adoption & understanding.

Why did you predict 42 for this data point?

*awkward silence*

Intrinsic ⟹ models are simple, so we interpret directly.

Post-hoc ⟹ train any "black box" model and analyze predictions

---

Local ⟹ why was a certain prediction made?

Global ⟹ how does the model work?

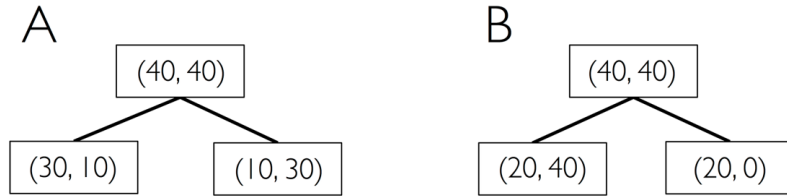$$y \approx \hat{y} = f(X) = w_0 x_0 + w_1 x_1 + \cdots \qquad w = w: \min_{w} ||\hat{y} - y||_2^2$$

Features with largest "effect"

are most important

effect $\equiv |w_i x_i|$

A

(40, 40)

(30, 10)    (10, 30)

B

(40, 40)

(20, 40)    (20, 0)

Features that ↓ impurity
the most are most important

Model-agnostic methods

$\hookrightarrow$ doesn't matter model type

Some model $\hat{y} = f(\vec{x})$

$\gtrless$

interpret $f$

# One global method: permutation importances

Idea: how much worse does my model get if a feature is "shuffled"?

1) Train a model, $\hat{y} = f(\vec{X})$

2) Compute error, $e_0 = \|y - \hat{y}\|_2^2$

3) For feature, $x_j$ in $\vec{X}$,
    a) Shuffle $x_j$
    b) Compute error, $e_j$ using $f(\vec{X})$

4) Importance $= e_j - e_0$

Fisher, Aaron, Cynthia Rudin, and Francesca Dominici. "All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously." http://arxiv.org/abs/1801.01489 (2018)

# One global method: permutation importances

## Compute on training data or validation data?

training $\Rightarrow$ how much is $X^j$ leveraged during training?

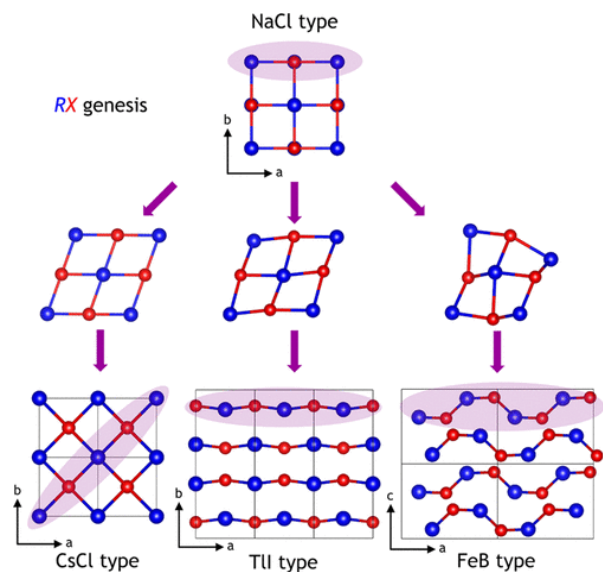validation $\Rightarrow$ how important is $X^j$ for generalization to new data?

Fisher, Aaron, Cynthia Rudin, and Francesca Dominici. "All models are wrong, but many are useful: Learning a variable's importance by studying an entire class of prediction models simultaneously." http://arxiv.org/abs/1801.01489 (2018)

# One global method: permutation importances

# From global to local methods

https://en.wikipedia.org/wiki/Simpson's_paradox

# Our data may not be globally interpretable

LIME $\Rightarrow$ Local Interpretable Model-Agnostic Explanation

Given: $\hat{y} = f(\vec{X})$

1) Select probe point, $\vec{X}_i$

2) Randomly generate $M$ new points

3) Predict $\hat{y}$ for these $M$ points

4) Train linear model on $\hat{y}_M$
   weighting $\mathcal{L}$ by proximity to $\vec{X}_i$

5) Interpret this simple model.

# Pitfalls of interpretable ML

## 1. Assuming one method will always work

Molnar et al., 2020, arXiv:2007.04131

# Pitfalls of interpretable ML

## 2. Bad model generalization

## 3. Unnecessary complexity

**Recall why we care about interpretability:**

Reliability → small change in $x_i$ shouldn't lead to a large change in $\hat{y}_i$

Causality → as we change $x_i$, can we anticipate change in $\hat{y}_i$

Trust ⟹ adoption & understanding

Why did you predict 42 for this data point?

*awkward silence*

**Intrinsic interpretability is always preferred!**