

Ensemble Machine Learning

- ↳ learning how to combine models to create stronger learners
- ↳ Perform strongly w/o overfitting

Outline:

- Bias Variance trade off
- Bootstrapping
- Bagging
- Rand. forests
- Ada Boost (Xg boost) (Cat boost)

Bias Variance Trade off

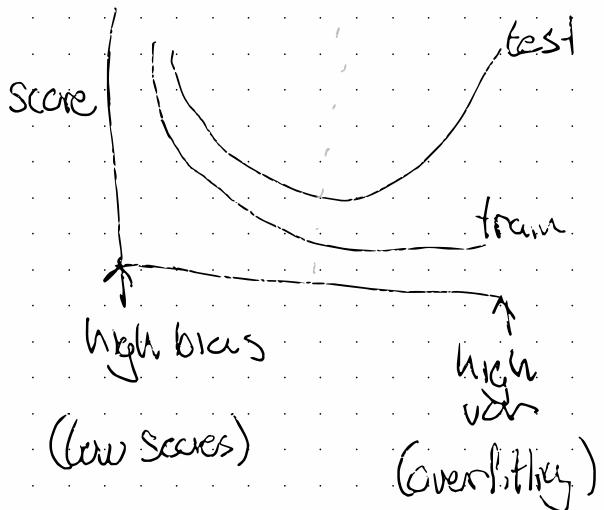
→ Bias	→ Variance	→ Irreducible Error
↳ error in model	↳ flexibility in model	↳ random error that cannot be accounted for
High biases = simple models	High var = complex models	

Bias variance tradeoff

Optimum

As models get more complex they don't generalise as well

→ goal, introduce just enough complexity for best score w/o overfitting



Irreducible error which is normally dist.

$$y = f(\omega_c) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

↑
data ↑
hypothetical
func

$$\hat{f}(x) = \text{model}$$

$$\therefore \text{err} = E[(y - \hat{f}(x))^2]$$

$$\therefore = E[(f(\omega_c) + \varepsilon + \hat{f}(x))^2]$$

Let $\hat{f}(x) = E[f(x)]$

then $\text{err} = E[(f(x) + \varepsilon - \hat{f}(x)) + \underbrace{\hat{f}(x) - f(x)}_{\text{Add & subtract to get the same thing}}]^2$

Add & subtract to
get the same thing

next, group so that

$$E[(f(x) - \hat{f}(x)) - (\hat{f}(x) - f(x)) + \varepsilon]^2$$

↪ expand & use properties of ε to
arrive at...

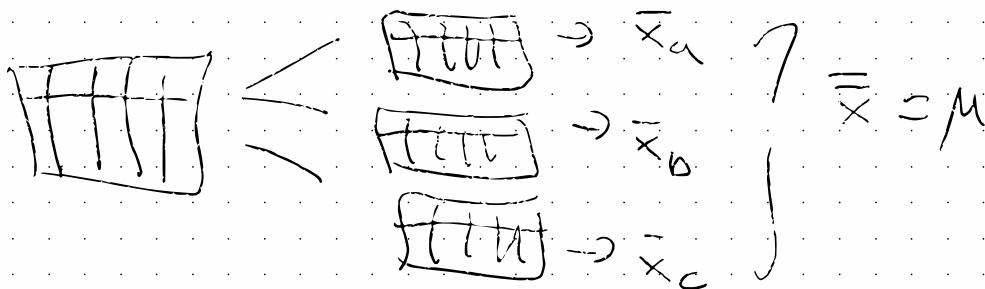
$$= \text{bias}^2 + \text{var} + \sigma_\varepsilon^2$$

How to find ideal parameters &/or
models?

↪ cross validation

Bootstrapping (Resampling)

- ↳ Possible to reduce both bias AND variance
- ↳ Bootstrapping = tool used to achieve this
 - ↳ Sample sub-datasets from original w/ replacement
 - ↳ possible to compute parameters from bootstrapped subsamples



→ Parameters:

$$E(\bar{\theta}_B) = \theta$$

$$\text{var}(\bar{\theta}_B) = \frac{1-p}{B} \sigma^2 + p \sigma^2 \quad \left. \begin{array}{l} p = \text{corr}(\hat{\theta}_i, \hat{\theta}_j) \\ \text{var}(\hat{\theta}_i) = \sigma^2 \end{array} \right\}$$

If $p=0$ then $\text{var}(\bar{\theta}_B) = \frac{1}{B} \sigma^2$ where

→ Derivation of variance

$$E(\bar{\Theta}_B) = E\left(\frac{1}{B} \sum_{i=1}^B \hat{\Theta}_i\right) = \frac{1}{B} E(\hat{\Theta}) = \Theta$$

suppose $E(\hat{\Theta}) = \mu$, $\text{var}(\hat{\Theta}) = E((\hat{\Theta} - \mu)^2) = \sigma^2$

& $\rho(\rho)$ = correlation between bootstrapped samples

($S_B = \sum \hat{\Theta}_i$, for notation's sake)

$$\text{var}(\bar{\Theta}_B) = E((\frac{1}{B} S_B - \mu)^2)$$

→ pull out $\frac{1}{B}$

→ expand sq. term

$$\text{var}(\bar{\Theta}_B) = \frac{1}{B^2} \in \left[S_B^2 - \underbrace{2\mu BS_B}_{\text{const.}} + \underbrace{\mu^2 B^2}_{\text{const.}} \right]$$

$$= \frac{1}{B^2} E[S_B^2] - \mu^2$$

$$E[S_B^2] = BE[\hat{\Theta}_i^2] + B(B-1) E_{i,j} [\hat{\Theta}_i \hat{\Theta}_j]$$

using equations defined before

$$E_{\hat{\Theta}_B} [\hat{\Theta}_B \hat{\Theta}_B] = \rho \sigma^2 + \mu^2$$

↳ misfit & melt

$$\text{var}(\hat{\Theta}_B) = \frac{1-\rho}{B} \sigma^2 + \rho \sigma^2$$

↳ Perfect corr

$$\text{then } \text{var}(\hat{\Theta}_B) = \sigma^2$$

when using decision trees, decision boundaries
don't correlate that well

↳ ends up reducing the variance

Bagging - "Bootstrap Aggregation"

↳ train multiple models & aggregate parameters (mean, etc.)

Pseudocode

models = []

for b in [1...13]:

model = model()

$\mathbf{x}_b, \mathbf{y}_b = \text{resample}(\mathbf{x})$

model.fit($\mathbf{x}_b, \mathbf{y}_b$)

models.append(model)

predict

R $\mathbf{y}_{\text{pred}} = \text{np.mean}$ (

[mdl.predict(\mathbf{x}) for mdl in models],
axis=1

Stacking

↳ Aggregates so far are not weighted
add weighting

$$f(x) = \sum_{m=1}^M w_m \hat{f}_m(x)$$

→ weight based on accuracy

$$\hat{w} = \arg \min E_{\text{pop}} [(\gamma - f(x))^2]$$

$$\hat{w} = E_{\text{pop}} \underbrace{[(F(x)^T F(x))]^{-1} F(x)^T \gamma}_{\text{looks like linear regression}}$$

looks like
linear regression

problem, we don't have $F(x)$

↳ Solution: (can't use $f(x)$ since only best will be selected)

train on all points except x_i, x_j

$$\hat{w}_{\text{stack}} = \arg \min \sum_{i=1}^n \left[y_i - \sum_{m=1}^M w_m \hat{f}_m(x_i) \right]^2$$

how do you then solve this equation

↳ quadratic programming problem

↳ don't worry about solving it.

→ recall: this method is similar to CGCG

Random Forest

- ↳ reduce correlation between different samples of Ob
- ↳ trees are good at capturing complex models w/ high depth, use this to create uncorrelated samples
- don't just select random rows, also select random features



Train random forest pseudo code

for b in $1 \dots B$:

$X_b, y_b = \text{resample}(X, Y)$

model = $\text{fittree}(.)$

while not at terminal node

and not max depth

- select random features

- choose best split from d features

- add split to model

append to models

Notes: Similar to bagging, but we also bag features

This method != bagged dt

↳ we've changed how splits are made

→ Method to build this class is very similar to original algorithm but splits change ↴

What happens if some columns are noise?

↳ This is why we use bootstrapping

Random Forest advantage

↳ Requires little tuning (Plug & play)

↳ learners can go to arbitrary depth w/o much penalty

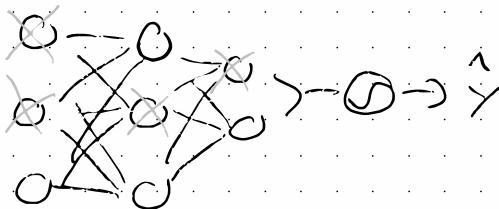
Good starting point above linear

class./reg. Z provides good info about
features. [↑] Start

"RF doesn't overfit" ← bagging/bootstrapping powerfully prevents overfitting

Dropout: Connection w/ DL

- ↳ allows for uniform training over a network so as to prevent overreliance on a specific node



- ↳ The way how we train via bootstrapping of features & observations is

AdaBoost

→ Also good "plug & play" Model

Bagging

↔

Boosting

low bias
high variance
weak learners

high bias
low var.
strong learners

Even after ensemble,
var will remain low

Hypothesis: comb of weak learners makes
a strong one

→ we weight each learner

$$F(x) = \sum_{m=1}^M \alpha_m f_m(x)$$

DT w/ max depth = 1

→ Advantage: trains really fast

Notable Differences:

- Labels = {-1, 1?} instead of {0, 1}
- $f_M(x) = \text{sign}\left(\sum_{m=1}^M \alpha_m f_m(x)\right)$
- Decision Boundary = 0

Auger Alm:

- Add base model 1 at a time
 - ↳ train on all data
 - ↳ weight each obs
 - ↳ modify weight on each rounds
 - ↳ larger weights for wrong classifications
 - ↳ model weight is a function of error

Additive Modelling

Adaboost loss function: Exponential Loss

classifications $\in \{-1, +1\}$

↪ cross entropy won't work, need smth else

$$L(y, f(x)) = e^{-yf(x)}$$

↑ same sign $\rightarrow 0$ never
opp sign $\rightarrow \text{inf}$] or 1

ie, we care abt sign more than how close it is to the real value

→ Application

$$(\alpha_m^*, f_m^*) = \underset{\alpha_m, f_m}{\operatorname{arg\min}} \sum_{i=1}^M e^{-y_i [f_m(x_i) + \alpha_m f_m(x_i)]}$$

↑ with model

Note that loss can be expanded due to exp

$$\begin{aligned} J &= \sum_{i=1}^N e^{-y_i f_{\text{true}}(x_i)} e^{-y_i f_{\text{true}}(x_i) \alpha_m} \\ &= \sum_{j=1}^N w_j^{(m)} e^{-y_j f_{\text{true}}(x_j)} \underbrace{\quad}_{\substack{\text{const} \\ \text{historical vals}}} \underbrace{\quad}_{\substack{\text{back to} \\ \text{root}}} \\ &\qquad \qquad \qquad y_i f(x_i) \text{ will always be} \\ &\qquad \qquad \qquad +1 \text{ or } -1 \end{aligned}$$

∴ separate into right & wrong sections

$$\begin{aligned} J &= e^{-\alpha_m} \sum w_j^{(m)} + e^{\alpha_m} \sum w_j^{(m)} \\ &\qquad \qquad \qquad y_j = f_{\text{true}}(x_j) \quad y_j \neq f_{\text{true}}(x_j) \\ &\qquad \qquad \qquad \downarrow \\ &\approx J = e^{-\alpha} A + e^{\alpha} B \end{aligned}$$

→ then, diff. & set to 0

$$\frac{\partial J}{\partial \alpha} = -e^{-\alpha} A + e^{\alpha} B = 0$$

$$\therefore \alpha_m = \frac{1}{2} \ln \left(\frac{A}{B} \right)$$

← weighted
correct
← weighted
incorrect

Comparison to Stacking

Stacking

- optimised using quadratic programming because of quad. loss func.

CUCV

- Not scalable
 $O(NM)$

Adaboost

- can be diff easier

- trained w/ sample weights

- Gradient algorithm
 $O(M)$

↳ lower complexity

Comparison to DNN

activation function \rightarrow tanh

$$\hookrightarrow -1 < \text{range} < 1$$

- each node in network is the output of a bayt algorithm
logistic reg

$$\text{Adaboost: } \hat{y} = \text{sign} \left(\sum_{m=1}^M \alpha_m \text{sign}(w_m^T x_c) \right)$$

$$\text{DL: } \hat{y} = \text{sign} \left(\sum_{m=1}^M \alpha_m \tanh(w_m^T x_c) \right)$$

→ hard vs soft outputs