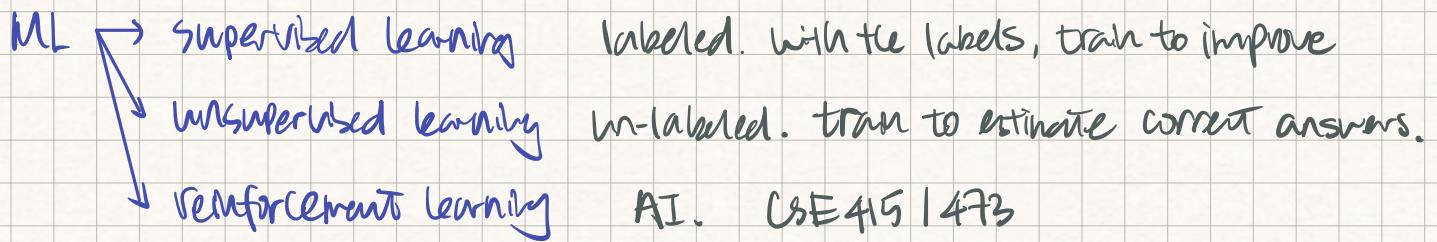


Learning Reflection 1

2022. 3. 28 (Mon)



1. Regression - Predict house prices
2. Classification - Pos. / Neg. Reviews.
3. Document Retrieval + Clustering - Find similar news articles.
4. Recommender Systems - Given past purchase, give recommend.
5. Deep learning - Image Recognition.

Models

Linear Regression,

regularized approaches (ridge, LASSO)

Linear Classifiers : Logistic regression

Non-Linear models : decision trees

Nearest neighbors, clustering

Recommender systems.

DL.

Algo.

Gradient Descent

Boosting

K-means.

Concepts

Point Estimation

MLE.

Loss Function,

bias-variance trade off,

Cross-validation.

Sparcity

Overfitting / Underfitting

Model selection.

Decision boundaries.

* Case Study 1 : Regression : Housing Prices

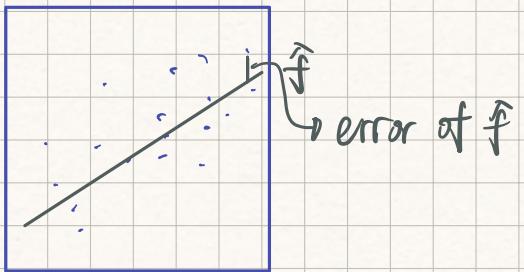
Goal: predict how much my house is worth.

Assumption: There is a relationship b/w YEAR & VAL , $y \approx f(x)$

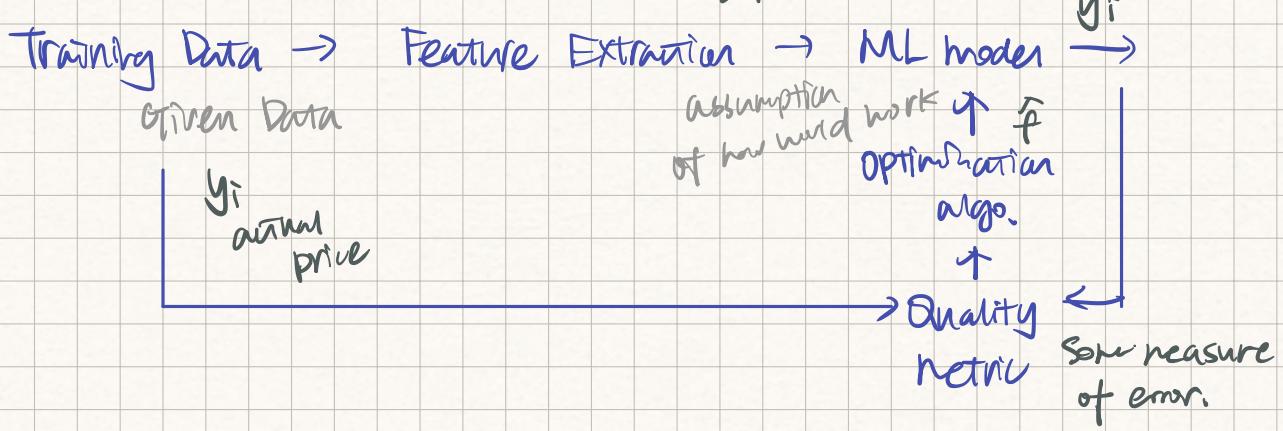
Model: $y \approx f(x)$ allows for some error of not an exact fit to true function.

Predictor: f is unknown. Predictor: \hat{f} , we predict: $\hat{y} = \hat{f}(x)$

Small error \rightarrow good fit.



* ML Pipeline:



* Regression.

Supervised learning algorithm.

Build predictor function learns how to map x^i to y^i

Linear Regression Model.

$$\hat{y}^i = w_1 x^i + w_0$$

Price Area of a house

$w_0 \rightarrow$ bias, starting price of a house

$w_1 \rightarrow$ slope/weight

$$\hat{y} = \hat{w}_1 x + \hat{w}_0$$

Cost / Loss of predictor - minimize the cost

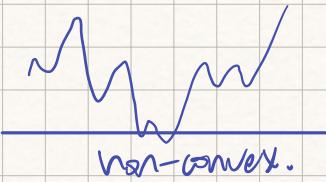
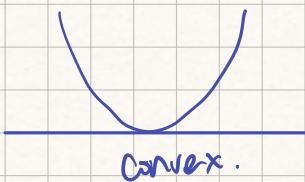
Default loss function for regression: MSE (Mean-Square Errors)

$$MSE = \frac{1}{n} \sum_{i=1}^n \underbrace{(y_i - \hat{y}_i)^2}_{\text{residual}^2}$$

y - actual val.

\hat{y} - predicted val.

Gradient Descent.



Start at some (random) weights

W While we haven't converged:

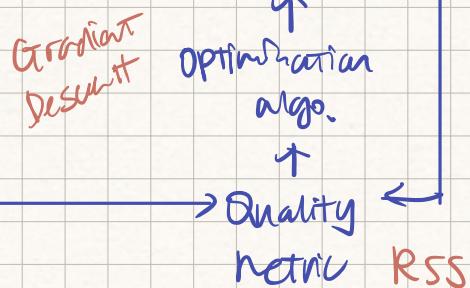
$$w = \alpha \nabla L(w)$$

$\alpha \rightarrow$ learning rate, controls how far you step.

$\nabla L(w) \rightarrow$ the gradients of loss function L on a set of weights w .

Pre-lecures.

Training Data \rightarrow Feature Extraction \rightarrow ML model \rightarrow linear regression.



* Polynomial Regression

$$\text{Model: } y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots + w_p x_i^p + \epsilon_i$$

use gradient descent to train.

* Features

$$\begin{aligned} \text{Model: } y_i &= w_0 h_0(x_i) + w_1 h_1(x_i) + \dots + w_p h_p(x_i) + \epsilon_i \\ &= \sum_{j=0}^p w_j h_j(x_i) + \epsilon_i \end{aligned}$$

$$\text{e.g.: } y_i = w_0 + w_1(\text{sq. ft.}) + w_2(\# \text{ bathrooms}) + \epsilon_i$$

Coefficients tell no. rate of change if all other features are constant.

• Distortion: is the difference between a data input & a feature

Data Inputs: columns of the raw data.

Features: values (possibly transformed) for the model
(done after feature extraction step).

Input: $\mathbf{x}_i = [x_{i1}, x_{i2}, \dots, x_{id}]$

Output: y_i

\mathbf{x}_i is the i th row

x_{ij} is the i th row's j th data input

$h_j(x_i)$ is the j th feature of the i th row.

• Recap . Linear Regression

Dataset: $[(x_i, y_i)]_{i=1}^n$ where $x \in \mathbb{R}^d$, $y \in \mathbb{R}$

Feature Extraction:

$$h(x) : \mathbb{R}^d \rightarrow \mathbb{R}^D$$

$$h(x) = (h_1(x), h_2(x), \dots, h_D(x))$$

Predictor

$$\hat{w} = \min_w RSS(w)$$

Regression Model:

$$\begin{aligned} y &= f(x) + \epsilon \\ &= \sum_{j=0}^D w_j h_j(x) + \epsilon \\ &= \mathbf{w}^T h(x) + \epsilon \end{aligned}$$

Prediction

$$\hat{y} = \hat{w}^T h(x)$$

Quality Metric

$$RSS(w) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2$$

ML Algo.

Optimized using Gradient Descent

☒ Assessing Performance

Why train ML models? do well on future data.

Minimizing RSS ≠ the model fits future data better

General loss function $L(y, \hat{f}(x)) = (y - \hat{f}(x))^2$

True Error: $E[L(y, \hat{f}(x))] = \sum_{x \in X} \sum_{y \in Y} L(y, \hat{f}(x)) p(x, y)$

→ all possible (x, y) pairs.

Need future data to calculate True Error, however we don't have future data in real cases. Try to minimize the true error is our goal.

Split dataset into train set & test set



Test Error:

$$RSS_{\text{test}}(w) = \sum_{i \in \text{Test}} (y_i - \hat{f}_w(x_i))^2$$

If the Test Set is large enough, this can approx. the true error.

☒ Train / Test Split.

Bigger test \Rightarrow Better estimate of true error.

But reduce the training set. Won't learn as well.

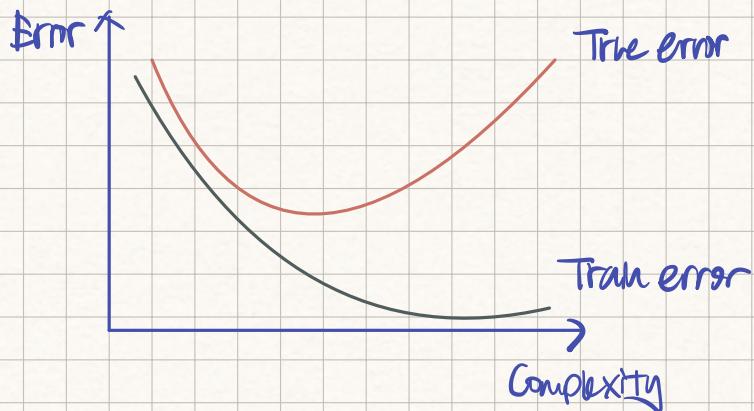
either $\begin{cases} 80: 20 \\ 90: 10 \end{cases}$

Never train model on data in test set.

2022/3/30 (Wed.)

Model Complexity . A model w/ more para. usually is more complex.

Train / True Error :



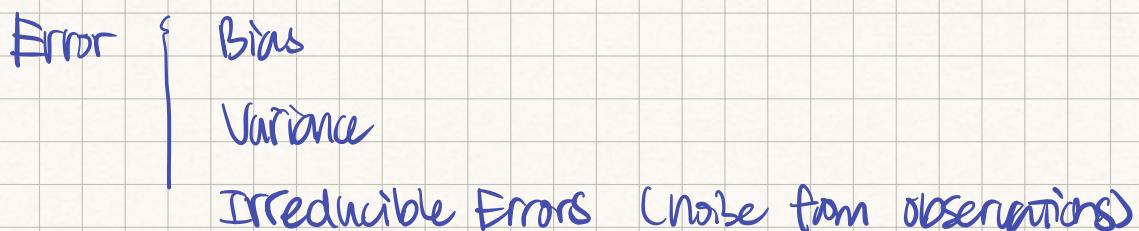
Underfitting : fix → remove uninformative features ;

increase complexity

(adding more features or introduce more non-linearities to capture non-linear patterns)

More complex → easy to overfit ; less complex → easier to underfit.

Bias - Variance Tradeoff.



* { Model too simple (low complexity), high bias. low variance
Model too complicated (high complexity), high variance. low bias

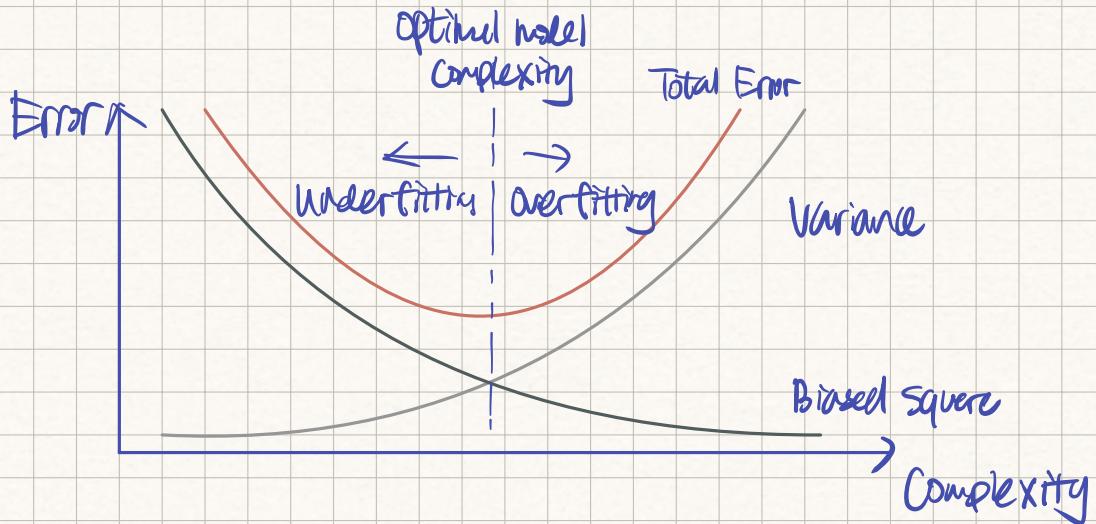
Bias : difference between the average prediction of our model and the expected value which we are trying to predict.

Variance: Variability in model prediction, meaning how much the predictions will change if a different training dataset is used.

Sources of errors for a particular model \hat{f} using MSE loss function:

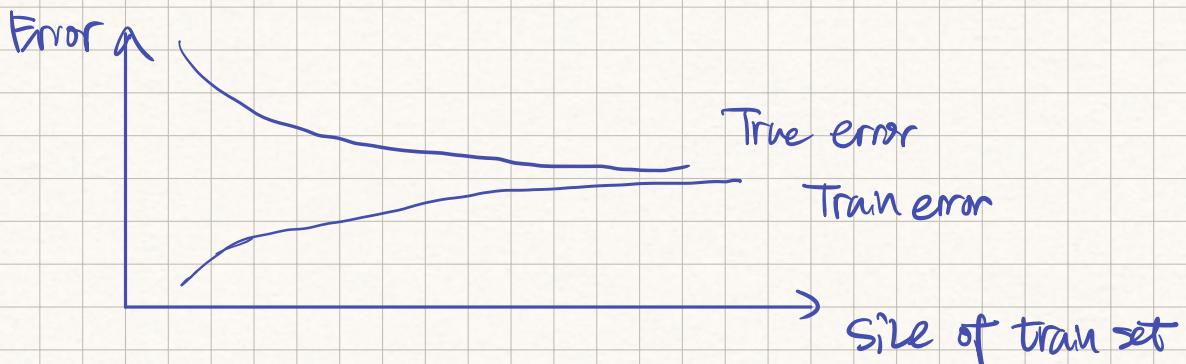
$$E[(y - \hat{f}(x))^2] = \text{bias}[\hat{f}(x)]^2 + \text{Var}(\hat{f}(x)) + \sigma^2_e$$

Error = Biased squared + Variance + Irreducible Error



* Model complexity doesn't depend on the size of the training set.

The larger the training set,
the lower the variance of the model. thus less overfitting .



* Choosing Complexity

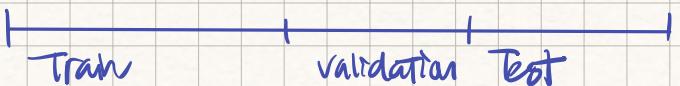
two ways to pick the model complexity w/o running our test set.

Using a validation set. \rightarrow Doing cross validation.

Validation set.



We can't use Test to choose our model complexity, instead,



↳ still fit a fully unbiased estimator of the true error, but less biased than the train set.

Pros: Easy to describe and implement

Pretty fast.

Cons: Have to sacrifice more training data.

Prone to overfitting

Summary

1. Use ML to build models and train them to well predict unseen (future) data.
2. Use True Error to make sure our models predict more accurately.
3. However, we need infinite data to calculate True Error but we don't have infinite data. Then, we need to estimate the True Error.

Therefore, we split the dataset into Train Set & Test Set with usually 90:10 or 80:20. Even though, a large Test Set decreases the error of estimating the True Error (can be more confident in the prediction), it decreases the accuracy of the Training Set.

4. Never use test error as the standard to pick models. We choose models only based on the Train Set. Then we use the Test Set for only testing ^{the} model.