

Day-14, Oct-16, 2024

Data Splitting: The Simplest way to Prevent the model overfitting and underfitting and to choose a model with the right level of complexity is to use data Spilthing.

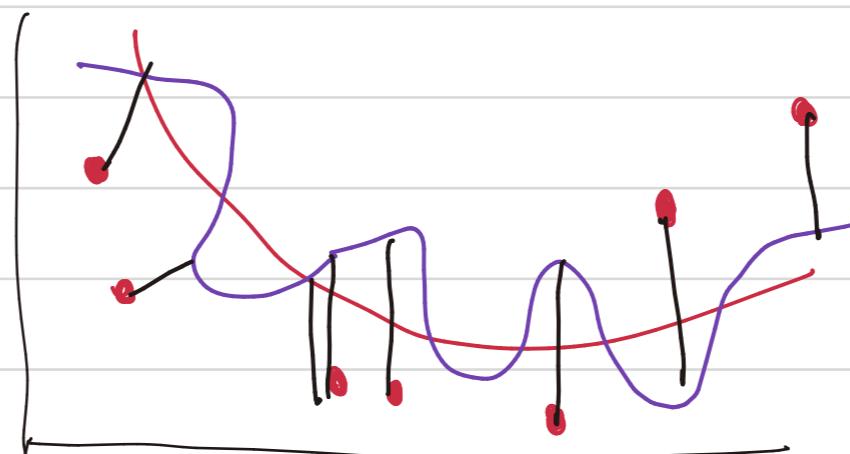
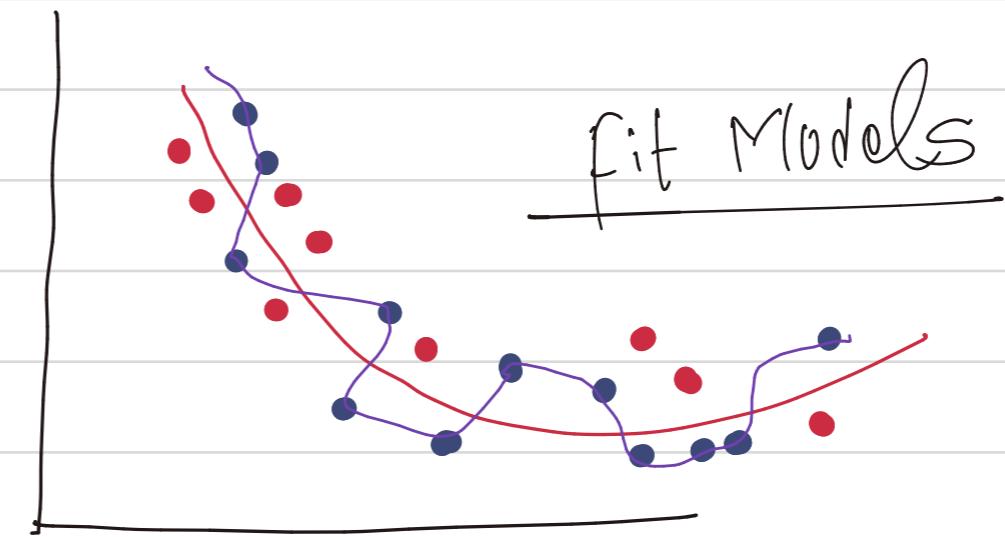
Randomly Split the data that are available into a train and a test set.

Fit a collection of models, all of different complexities, on the training set

look at the performance of those models on the test set

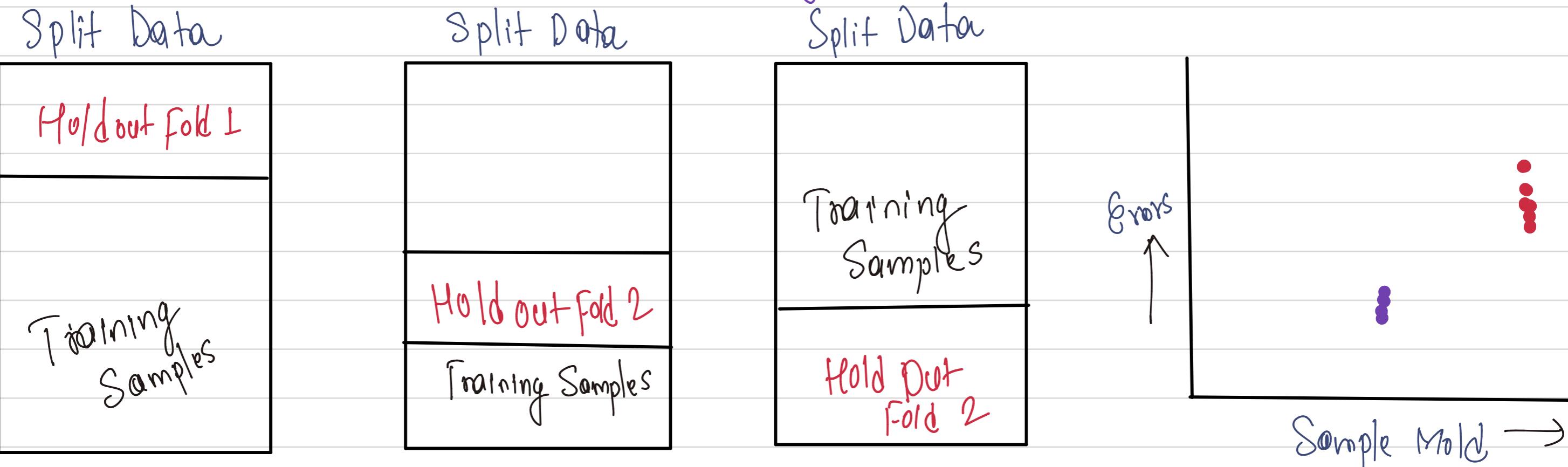
choose the model with the best performance on the test set.

Split Data.



We only train the training Samples

A more Sophisticated approach is to use cross-validation. Instead of using train/test Split. We can group the data into K-different "folds". For each Complexity level, we can train K each time leaving out one of the folds. Performance of these models on the estimating the appropriate Complexity level.



How does Cross Validation Work?

1) Data Splits into k -folds (K -Subsets)

2) Training and Testing:

Model is trained on $k-1$ folds (eg: 4 out of 5 folds)

and (tested on the remaining 1 fold)

This process is repeated K times, each time using a different fold

as the test set and the remaining $k-1$ folds for training

3) Average Results:

After all K -iterations, the model's performance (eg: accuracy)

error) is averaged across all folds to give a more reliable estimate of its generalization ability.

Theoretical part:

There are a few subtleties related to train/test splits and cross-validation that are worth knowing -

- o We want the held-out data to be as representative of new, real-world data as possible.

Sometimes, randomly split data will not be representative - for example, maybe the data are not collected over time.

A more "honest" split would be into past vs. future samples

starting at a few different timepoints; training on past and evaluating on future. The sample issue occurs if we plan on applying a model to a new geographic context or market segment, for example.

Cross-validation can be impractical on larger datasets. For this reason, we often see only a single train/test split for evaluation of algorithms in large data settings.

We incur a bias when splitting the data. The fact that we split the data means that we aren't using all the available data, which leads to slightly lower performance.

if it hurts more complex models more than simpler ones, than it can even lead to slightly incorrect selection.

After having chosen a given model complexity, it's worth retraining that model on the full available dataset.

Randomly Split the data before training the model:

Example of future and past Splitting of dataset is on predicting the Stock price, Geographic data-set,

Be-careful with the data splitting like using subset of datasets, how and what to use and where to use dataset?

After the Coding part for both Cross-Validation and train-Split/Test

Note!:

Related to overfitting, there is important phenomenon called the bias-variance trade off.

To understand this bias-variance trade off, imagine being able to sample hundreds of datasets similar to the one that is available. On each of these datasets we can train a new model (at a given complexity level). Then, we can define

① Bias

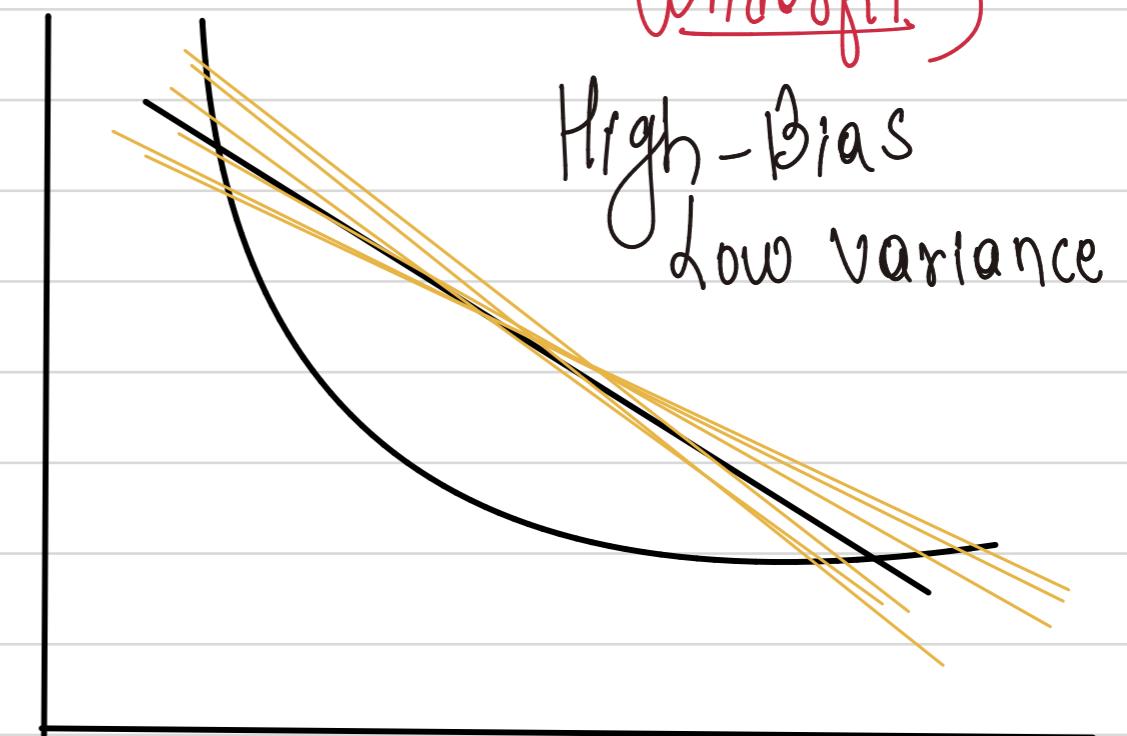
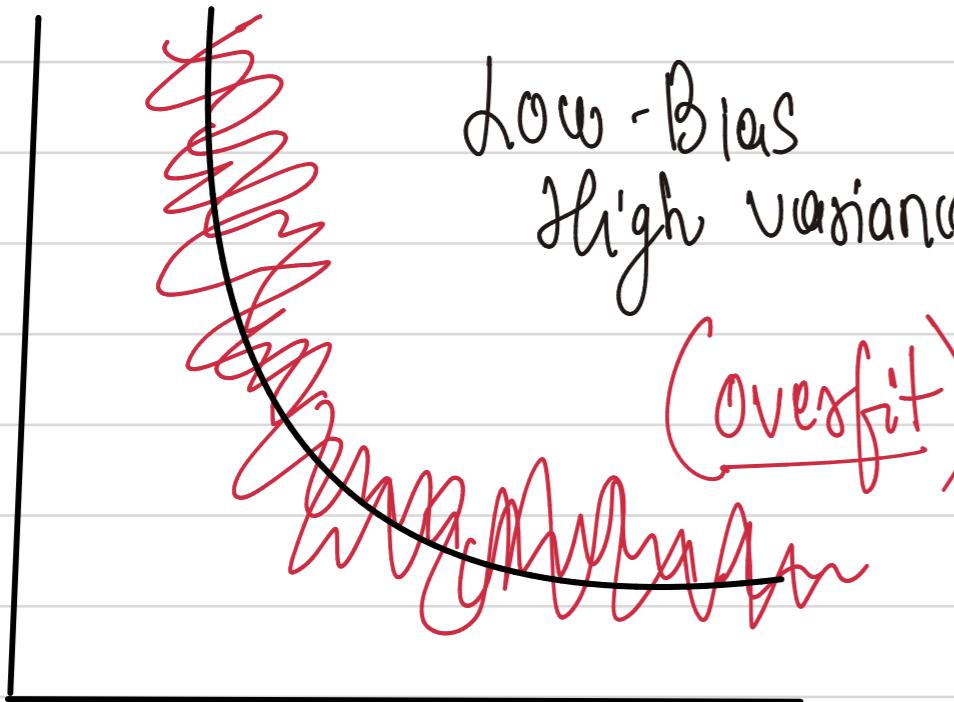
② Variance

① Bias:

Average the predictions from models trained across each of the imaginary datasets. How far off they are from the truth, on average?

② Variance:

How similar are the predictions from across the different runs?



Ideally, we would have low values for both bias and variance,

Since they both contribute to the performance on the test-set.

In practice though there is trade-off.

Often, the high-powered models that tend to be closest to the truth on average might be far off on any individual run (high variance). Conversely, overly simple models that are stable from the run to run might be consistently incorrect in certain regions (bias).

Models with high variance but low bias tend to be overfit &
Models with low variance but high bias tend to be underfit.

H Models that have good test or Cross-Validation errors have found
a good compromise between bias and variance.

Tips: Starts with the Complicated or Complex Models
Shows to him and optimize this, then we could have more
Generalized Model. Rather than Starting from Simplex to the
Complex Model.

Choosing the Models with the right Complexity.

Logistic Regression Need & Loss Function

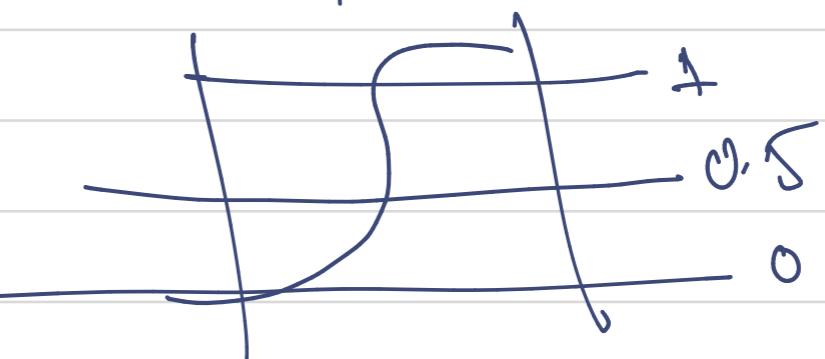
- ① There was no discontinuous function in the linear function because of its threshold functions.
- ② $h_w(x)$ is not differentiable this makes learning with perception very complicated.
- ③ We can use logistic function (or logit) and can soften the hard threshold functions - approx Hard function with continuous and differentiable functions.

So, the Soft threshold functions -

$$\text{Logistic } (\lambda) = \frac{1}{1+e^{-\lambda}}$$

Where $\lambda = wx$ ≠ weight * input

- With these our output is no longer 0 or 1 (linear function used before as perceptron learning with XOR and AND gates).
- Output is between 0 to 1 such that output can be 0.23, 0.5, 0.78, 0.9 so



$g \Rightarrow$ Logistic function

$g' \Rightarrow$ derivative of logistic regression

So to $h_w(x)$ Continuous we use loss function

Chain Rule $\Rightarrow \frac{\partial g}{\partial x} (f(x)) \Rightarrow g'(f(x)) \cdot \frac{\partial f(x)}{\partial x}$

Derivation of loss function

$$\frac{\partial}{\partial w_i} \text{loss}(w) \Rightarrow \frac{\partial}{\partial w_i} (y - h_w(x))^2$$

$$\frac{\partial}{\partial w_i} \text{loss}(w) \Rightarrow \frac{\partial}{\partial w_i} (y - h_w(x))^2 \quad \left[\frac{\partial z^2}{\partial z} = 2z \right]$$

$$\Rightarrow 2(y - h_w(x)) \times \frac{\partial}{\partial w_i} (y - h_w(x))$$

$$\Rightarrow -2(y - h_w(x)) \times g'(w \cdot x) \times \frac{\partial}{\partial w_i} w \cdot x$$

$$\Rightarrow -2(y - h_w(x)) \times g'(w \cdot x) \times x_i$$

derivative of 'g' $\rightarrow g'$

$$g \Rightarrow \frac{1}{1 + e^{-z}}$$

$$g'(z) \Rightarrow g(z)(1 - g(z))$$

So, we have,

$$\begin{aligned} g'(w \cdot x) &\Rightarrow g(w \cdot x)(1 - g(w \cdot x)) \\ &\Rightarrow h_w(x) \cdot (1 - h_w(x)) \end{aligned}$$

Weight update for minimizing the loss is -

$$w_i \leftarrow w_i + \alpha (y - h_w(x)) \times h_w(x) \times (1 - h_w(x)) \cdot x_i$$

$$\text{So, } w_i \leftarrow w_i + \alpha J$$

where $J \Rightarrow (y - h_w(x)) \times h_w(x) \times (1 - h_w(x)) \cdot x_i$

Logistic Regression are better for larger dataset where data are noisy and non-separable

$h_w(x)$

where

x is the input feature vector

$w \Rightarrow$ weight vector (model parameters)

$h_w(x)$ output (predicted probability) from the logistic

Regression model of given input (x)

Squared Error loss: $(y - y')^2$

where y' is the predicted outcome & y is the actual outcome

So,

$$\text{loss}(w) = (y - h_w(x))^2$$

y can be either '0' or '1' & $h_w(x)$ is the predicted value [0,1].

We calculate gradient

$$\frac{\partial}{\partial w_i} \text{loss}(\omega)$$

Expand the loss function:

$$\text{loss}(\omega) \Rightarrow (y - h_w(x))^2$$

So, we apply Chain Rule-

$$\frac{\partial^2}{\partial z^2} = \frac{\partial^2}{\partial z^2}$$

So we get,

$$\frac{\partial}{\partial w_i} (y - h_w(x))^2 \Rightarrow 2(y - h_w(x)) \cdot \frac{\partial}{\partial w_i} (y - h_w(x))$$

Because taking derivative, we take first after function derivative-

$$z^2 \Rightarrow 2z$$

$$(y - h_w(x))^2 \Rightarrow 2(y - h_w(x))$$

But we have inner function too $(y - h_w(x))$

Outer function $z^2 \Rightarrow dz$

$$z = (y - h_w(x)) \Rightarrow d(y - h_w(x))$$

Inner function: $\underline{z} \Rightarrow y - h_w(x)$

derivative with w_i -

$$\frac{\partial}{\partial w_i} (y - h_w(x)) \Rightarrow \underline{-h_w(x)} \frac{\partial}{\partial w_i}$$

$$S_0, \frac{\partial}{\partial w_i} (y - h_w(x))^2 \Rightarrow 2(y - h_w(x)) \cdot \frac{\partial}{\partial w_i} (y - h_w(x))$$

$$\Rightarrow 2(y - h_w(x)) \cdot -\frac{\partial h_w}{\partial w_i}$$

$$\Rightarrow -2(y - h_w(x)) \cdot \frac{\partial h_w}{\partial w_i}$$

$$\Leftarrow -2(y - h_w(x)) \cdot g'(w \cdot x) \cdot x_i$$

$\left[\frac{\partial}{\partial w_i} (y - h_w(x)) \right]$
 is $- \frac{\partial h_w}{\partial w_i}$

and $h_w(x) \Rightarrow g(w \cdot x)$.

$\frac{\partial h_w}{\partial w_i}$ is $g'(w \cdot x)$