

Generalizing to unseen data

Reducing Train Loss \rightarrow Done?

ML \rightarrow Generalize to unseen data

Overfitting & Underfitting



Train Test
 \hookrightarrow same dist \leftarrow ML's assumptions

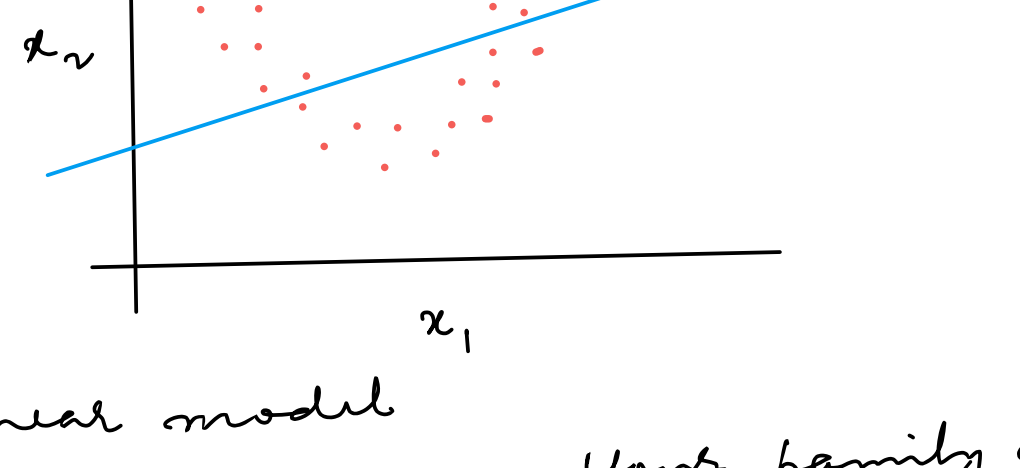
Very good \uparrow Test $\downarrow \leftarrow$ Overfitting

Poor on train \downarrow Test $\downarrow \leftarrow$ Underfitting

Poor train \downarrow Test \uparrow ? Lucky set

Bias & Variance

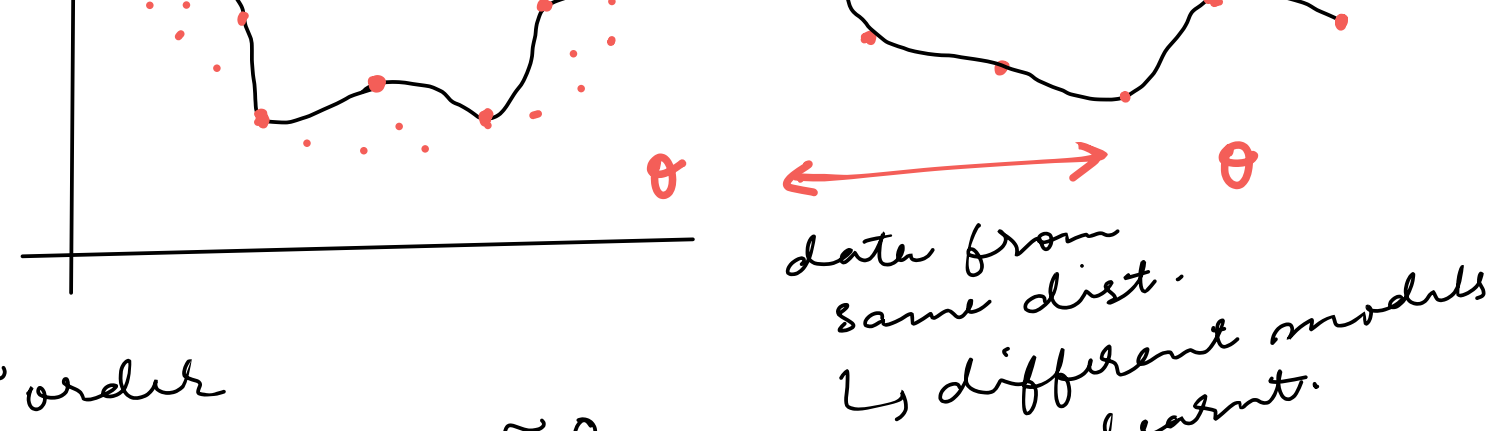
$$\theta_0 x_0 + \theta_1 x_1 \leftarrow \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_0^2 + \theta_3 x_1^2 + \theta_4 x_0 x_1$$



Linear model

\hookrightarrow high bias \rightarrow Your family of models fundamentally cannot approximate $f(x)$

$$y = h(x)$$



5th order

Training error ≈ 0

Training data \rightarrow Sample of an infinite set

Finite set \rightarrow Randomness

$$y = \underbrace{h(x)} + \epsilon_i \leftarrow \text{Fit the noise}$$

Bias \uparrow

Variance \rightarrow Increase training $h(x) \neq f$

Bias: Your test set error if you trained your model on virtually infinite data

Variance: Captures your model's sensitivity to randomness.

Model Validation / Selection

100% \leftarrow Train x

Train Validation **Test**

80% train 10% valid \leftarrow cal loss } cheating

10% test \leftarrow After many iterations on the previous two.

1. Hold out cross validation

come out a valid set.

quad on val set

Train \leftarrow Linear

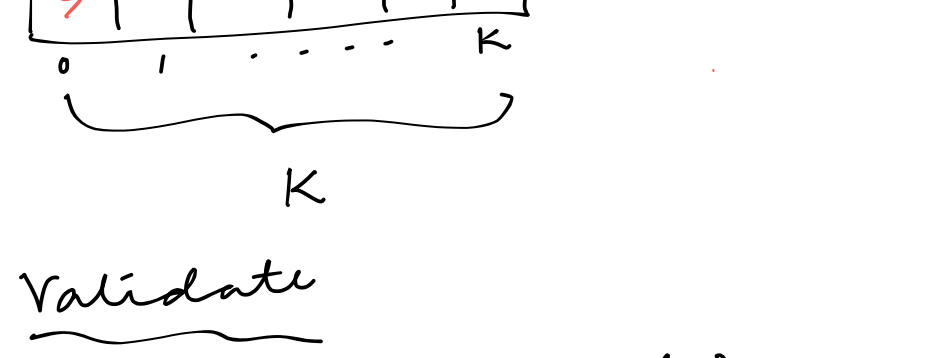
Val \leftarrow Quad

Val \leftarrow Linear beats Quad

80% 10% 10% }
 800 100 100 }
 1,000,000 }
 1,000,000 }
 1,000,000 }
 Valid Test
 Numbers \rightarrow 1000 data

2. K-fold cross validation

Hold out wastes data



Validate

0th block \rightarrow valid } val loss

rest \rightarrow train

1st block \rightarrow valid } val loss

rest \rightarrow train

average

3. Leave one out cross val

n data points

n-1 train } repeat n times

1 valid

If \rightarrow Model selection \rightarrow Linear / Quadratic

Linear \rightarrow Train + val

Retrain

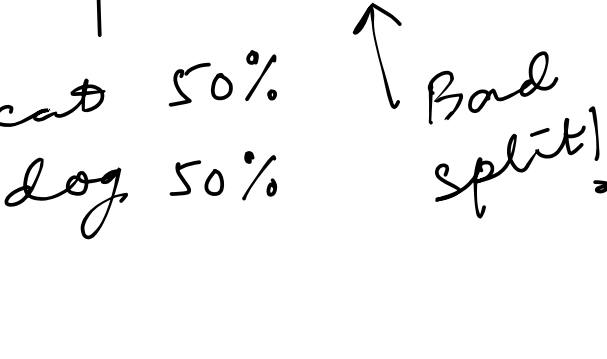
Deploy

Stratification

Respect the distribution

1. shuffle data

2. Stratification



Bad split!

80% \rightarrow dogs

20% \rightarrow cat

80% accuracy

Train set respects the dist.

Valid set skewed 90% cats \leftarrow 90%

60% 20% 20%

Train Val Test

60:20:20 60:20:20 60:20:20

F1 score

Cross val done \rightarrow Overfit

Train \checkmark

Val \downarrow

1. Get more data.

2. Regularization \leftarrow Restrict your model's power.

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2 + \lambda R(\theta)$$

$R(\theta) \rightarrow$ Regular. \uparrow

simple $\rightarrow R(\theta) \downarrow$

$$\text{Regulariser} = \|\theta\|_2 = \left(\sum_{i=0}^d \theta_i^2 \right)^{1/2} \leftarrow \text{Ridge regression}$$

$$\|\theta\|_1 = \sum \theta_i \leftarrow \text{Lasso regression}$$

$$J(\theta) := J(\theta) + \lambda \|\theta\|_1^2$$

θ_i should be high

Will $R(\theta)$ still push it down.

Dataset \rightarrow Extra features \leftarrow Feature selection

length

Breadth

$R(\theta) \rightarrow$ side effect \rightarrow Feature selection

\rightarrow Lasso $\rightarrow \|\theta\|_1 \rightarrow 0 \rightarrow 0$

$\|\theta_2\| \rightarrow \approx 0$

Why does Lasso do feature selection

& ridge regression doesn't.

1. Minimizing train loss isn't the final goal.

\hookrightarrow Generalize to unseen data.

2. Is my model overfitting or underfitting

\hookrightarrow Bias and variance

3. Bias \rightarrow your model is fundamentally incapable of representing the data.

\hookrightarrow Pump more data? X

4. Variance \rightarrow Sensitive to randomness in the finite data set.

\hookrightarrow Pump more train data. \checkmark

\hookrightarrow Regularize

\hookrightarrow Use a simpler model.

5. How do get an idea of bias/var?

\hookrightarrow Cross-validation

1. Hold out \rightarrow Single val.

2. K-fold \rightarrow K valid } average

3. Leave one out \rightarrow use 1 datum } repeat

at val set } for each datum

\hookrightarrow Always shuffle your data before splitting.

\hookrightarrow In case of classification \rightarrow Stratify

6. Regularization \rightarrow L_2 norm of $\theta \rightarrow$ Ridge regression

L_1 norm of $\theta \rightarrow$ Lasso regression

7. Lasso automatically does feature selection.

Why? \rightarrow Resource

Code LR \rightarrow sklearn

Linear regressor

LR \rightarrow Random θ

update rule $\rightarrow \theta$ improves

convergence criteria

Class

Numpy

Pandas