# ✦ Bias/Variance:-

$(\theta_0 + \theta_1 x_1)$

$(\theta_0 + \theta_1 x_1 + \theta_1 x_1^2 \cdots \theta_5 x_1^5)$

← $

Underfits
(High Bias).



$

Size

$

$(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2)$

Size
↓
Overfits
(High
Variance).

$



Size

(Just Right).

(*) The term high bias comes from the fact that our model was pr $\boxed{\text{Supposedly bias}}$ towards less complearty (linear) while actual data distribution was more complex (Quadratic). (High Bias).

(*) The term high bias Variance mean that everytime we fit our model on slight slightly different data. our $\boxed{\text{model would change}}$ to a large extent. Therefore our model fittiny and predicting is $\boxed{\text{highly varied}}$ given even a slight change in the distribution of data. (High Variance)

(*) $\$$

L·R (Linear) ⟶ Underfits.

⟶ Just right.



CL·RC Non-Linear High Polynomial Power ⟶ Overfits.

Size.

Fig

⑤ Regularization (Preventing Overfitting):-

(a) $\min \frac{1}{9} \frac{1}{a} \|h_\vartheta(x) - y\|^2$

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^{m} \| y^{(i)} - \theta^T \cdot x^{(i)} \|^2 + \underbrace{\frac{\lambda}{2} \| \theta \|^2}_{}$$

$$\underbrace{\phantom{\min_{\theta} \frac{1}{2} \sum_{i=1}^{m} \| y^{(i)} - \theta^T \cdot x^{(i)} \|^2}}_{\text{Cost}} \qquad \text{Regularization term.}$$



$\lambda = 0 \rightarrow$ Overfitting

$\lambda = 1 \rightarrow$ Right

$\lambda = 100 \rightarrow$ Underfitting.

↳ Overfitting preventing by decreasing the size of $\theta$.

↳ Very large $\lambda$ mean $\theta$ will be forced to be close to zero.

ⓧ $\max_{\theta} \sum_{i=1}^{m} \phi + \theta$

ⓧ $\arg\max_{\theta} \sum_{i=1}^{m} \log P(y^{(i)} | x^{(i)} ; \theta) - \lambda \| \theta \|^2$.

↳ SVM does not overfit badly because

$\min \| w \|^2$ objective function have same effective as regularization.
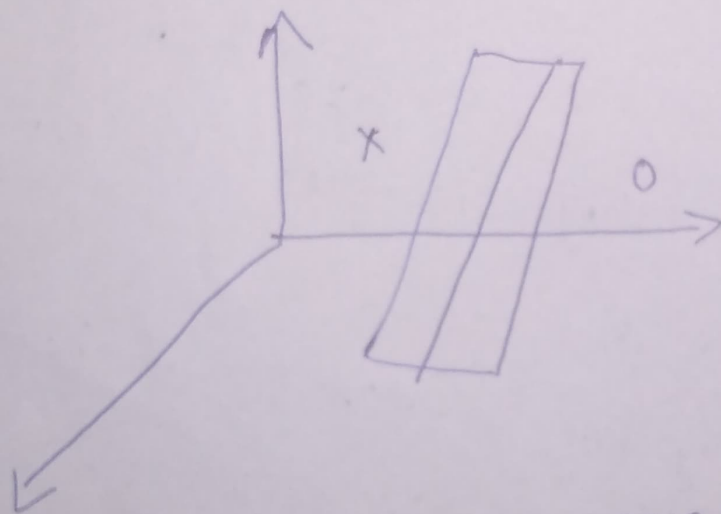
Ⓧ Text classification

⊙ m = 100

n = 10,000.

$X = \begin{bmatrix} 1 & 0 \\ 0 & adan \\ \vdots & \vdots \\ 1 & 1 \\ 0 & \\ \vdots & \end{bmatrix}$

→ if you fit logistic regression to such a data where n >> m then it will overfit the data ~~becoue~~

→ But if you use logistic regression with regularization then it is a good algorithm for text classification.

Ⓧ



Ⓨ More Example. mean logistic regression even without regularization won't overfit and would preform good provided m > n.

Ⓩ Features al different scales should he normalized on the same range so that these learning models trains faster and perform better.

→ Regularization and θ prior.

④ $S = \{(x^{(i)}, y^{(i)})\}_{i=1}^{m=1}$

$$P(\theta|S) = \frac{P(S|\theta)\, P(\theta)}{P(S)}$$

⊙ $\arg\max_\theta P(\theta|S) = \arg\max_\theta P(S|\theta) \cdot P(\theta)$

$$= \arg\max_\theta \underbrace{\left( \prod_{i=1}^{m} P(y^{(i)}|x^{(i)};\theta) \right)}_{\text{Generalized Linear model}} \cdot P(\theta).$$

$P(\theta): \theta \sim \mathcal{N}(0, \tau^2 I)$

$$P(\theta) = \frac{1}{\sqrt{2\pi}\,(\tau^2 I)^{1/2}} \exp\left( -\frac{1}{2} \theta^T (\tau^2 I)^{-1} \theta \right).$$

⊛ If this is the prior distribution for θ and you plug it in then it is equivalent to regularization.

⊛ Regularization is equivalent to assuming a prior θ Gaussian distribution.

FP

## Frequenist vs Bayesian Approach:

→ In Frequenist Approach we want to find the value of $\theta$ which makes that the data as likely as possible
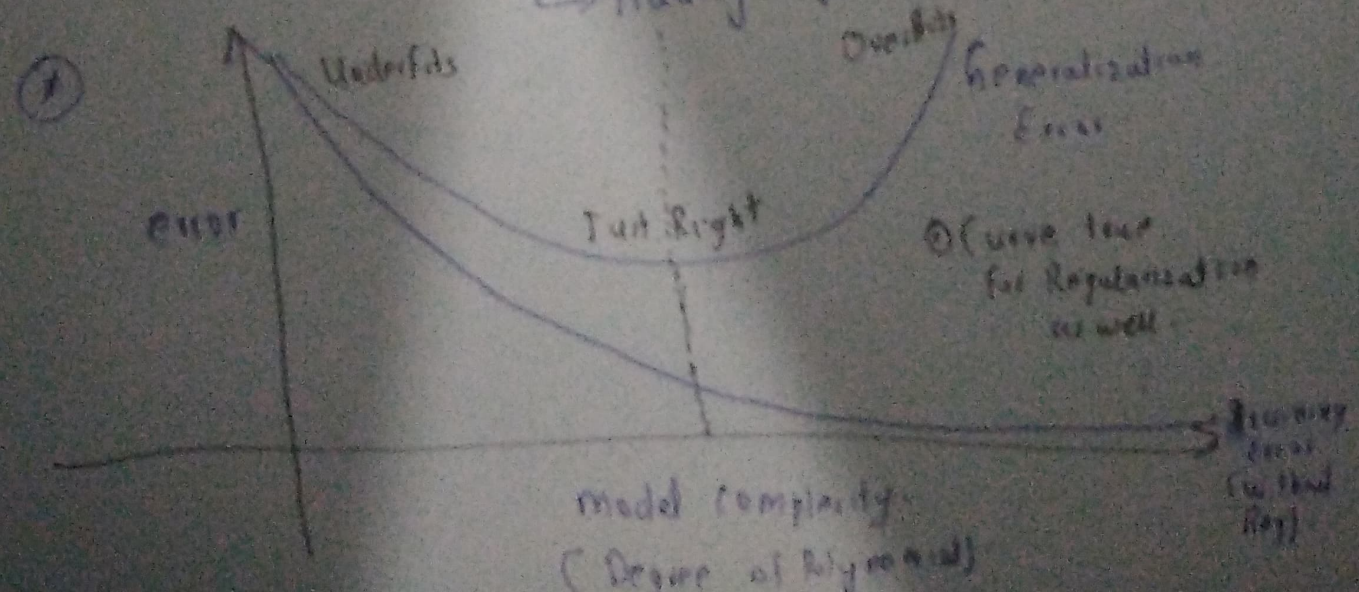
$$\underset{\theta}{argmax}\, P(S|\theta) - MLE \quad (Maximum\ Likelihood\ Estimate)$$

→ In Boysenia approach you already have some prior belief about $\theta$ in form of $P(\theta)$ (Gaussian prior).

↳ Our goal is to find

$$value\ \underset{\theta}{max}\, P(\theta|S) \to MAP\ (Maximum\ a\ Posteriori$$

↳ Adding Regularization



- error (y-axis)
- Underfits
- Overfits
- Generalization Error
- Just Right
- ⊙ Curve true for Regularization as well
- model complexity (Degree of Polynomial)
- training error (without Reg)

(a) **Traning / Dev / Test**

→ Spliting the Dataset into different
  Subset:

→ 10000 examples

Different ←
Hyperparameter
needed to
be choose
$$\begin{cases} Eg \quad \theta_0 + \theta_1 x \\ \quad \theta_0 + \theta_1 x + \theta_2 x^2 \\ or \quad choosing \quad value \quad \lambda \\ or \quad T \end{cases}$$

→ Splity your training data:
  ① $S_{train}$ · $S_{dev(development)}$ · $S_{test}$

② Train each model(option for
  different hyperparameters) on $S_{train}$.
  └ Get some hypothesis h.

③ Mesure     error   on $S_{dev}$ for
  new developed hypothesis h.

④ Pick h with lowest error on $S_{dev}$.

(*) If you want to check the
  unbiased peiformance of your model
  then evaluate on $S_{test}$ since $S_{dev}$ has
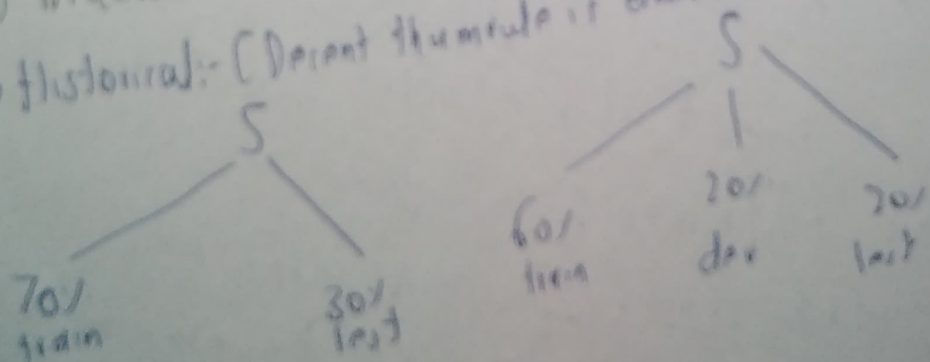  alredy been optimized upon by your model.

④ $S_{test}$ gives the real performance of your model. Since $S_{dev}$ is a biased estimate as model hyperparameter are adapted to perform on $S_{dev}$.
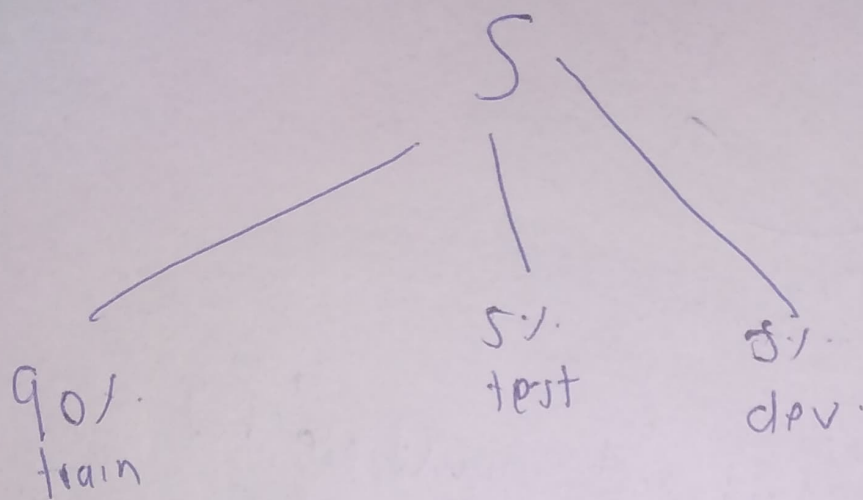
⑥

| degree | $S_{dev}$ Error | $S_{test}$ Error |
|--------|-----------------|------------------|
| 1 | 10 | 10 |
| 2 | 5.1 | 5 |
| 3 | 5.0 | 5.0 |
| 4 | (4.9) →Biased estimate | (5.0) → Actual Performance |
| 5 | ⚡7 | 7 |
| 6 | 10 | 10 |
| | ⋮ | ⋮ |

✗ Sometimes you can ignore $S_{test}$ in real world product but in publishing paper you need to have $S_{test}$.

✗ How much data to each split?
→ Historical: (Decent thumbrule if dataset is small e.g thousands examples).

ⓧ Era of Big-data (10 million example):-

S

90%.
train

5%.
test

5%.
dev·

ⓨ

ⓧ But, Sometimes you need a largest test
dataset if you the performance improvement
between models is really small.

└ ⓧ Choose test and dev big
enough to make meaningful
comparision between different models.

ⓧ This Process of holding dev set from
train set is called hold-out cross-
validation (Simple hold-out cross validation).
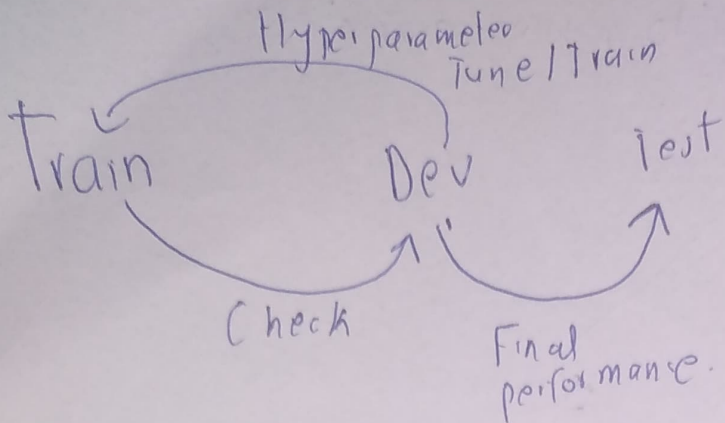└→ Development set = Cross validation
(Dev) Set

✗

Hyperparameter
Tune/Train

Train ⟶ Dev ⟶ Test

Check
Final performance.

Fig.

—Say Evaluate on test set multiple time but don't make decision or tune hyper-parameter based on the performance of test which is supposed to be an unbiased estimator benchmark

⟶ Dev set evaluation is used for hyper-parameter tunning

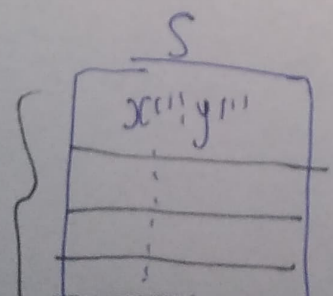✗ Small datasets:- (K-Fold - CV-(Cross-Validation).

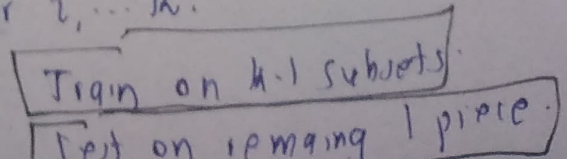① m=100 (Expensive Examples).

70% — $S_{train}$ , 30% — $S_{dev}$.

② This next procedure should be used only with small datasets:-

③ K-Fold CV (Cross-Validation
Divide/split data into five different subset

S

| $x^{(i)}, y^{(i)}$ |
|---|
| : |
| : |
| : |
| : |

For i, ... k:-

Train on k-1 subsets
Test on remaining 1 piece.

5-k subsets
of each
20 examples.

→ You will cross-validation for different hyperparameter and calculate average of metric for the cross-validations and then see which hyperparameter performed the best.

→ Final Optional Step:-

→ Refit the Model with chosen hyperparameter on Complete dataset.

→ This procedure make efficient use of data but is computationally very expensive.

Ⓧ Extreme version of EV K-fold CV :-
(For really small dataset):-

   m = 20 (20 examples)

   ⊙ Leave-one-out CV :-

   k = m → Divide your data into as many pieces as your training instances and leave only instance for evaluation at each iteration of validation.

⊛ **Feature Selection:-**

⊚ Small subset of features that are most useful for your task

⊚ Most Important

⊙ Depends on domain & context.

① **Foward Search.**

→ Start with $\mathcal{F} = \emptyset$

Repeat {

1) Try adding each feature $i$ to $\mathcal{F}$, which single feature addition most improve dev set performance

2) Add that feature to $\mathcal{F}$ }

**Example**

⊚ $X_1 \cdots X_5$ (Five features.

$\emptyset$     $h(x) = \theta_0$ (No feature)

$$
\begin{bmatrix}
\emptyset + x_1 \\
\emptyset + x_2 \longrightarrow h_2(x) = \theta + \theta_2 x_2 \quad \mathcal{F} = \{x_2\}. \\
\vdots \\
\emptyset + x_5
\end{bmatrix}
$$

$$
\begin{bmatrix}
x_2 \quad x_1 \longrightarrow \theta + \theta_1 x_1 + \theta_2 x_2 \quad \mathcal{F} = \{x_1, x_2\}. \\
x_2 \quad x_3 \\
x_2 \quad x_4 \\
\phantom{x_2} \quad x_5
\end{bmatrix}
$$

→ The above example shows that you start with empty array of features and then at each first iteration train five different model with only one feature respectively. See which model with one specific feature perfoim the best and add that feature to your feature array

→ Then at second iteration you add an other feature one at a time to previously choosen feature and see which model perform best with these two feature (one → new, one → previously choosen)

  ↳ Yoy continue until you have either added all features or performance jump by new features addition is minimal.

Ⓧ Backward → Start with all features and remove features to one by one to see the performance degradation.

$$J = (X_1, X_2, X_3, X_4, X_5).$$

Repeat {

  1) Remove Feature ith and train model with remaining features
  2) Remove feature with most minimal decrease in performance when removed }

Example

$$\mathcal{F} = (X_1, \dots X_5) \in (\text{Five Features}).$$

$$\begin{bmatrix} X_1, X_2, X_3, X_4 \\ X_2, X_3, X_4, X_5 \\ X_3, X_1, X_2, X_5 \\ X_1, X_2, X_4, X_5 \\ X_2, X_3, X_4, X_5 \end{bmatrix} \longrightarrow$$ Five different model and see whose accuracy is still the greatest and remove that feature (Say feature 5).

$$\mathcal{F} = (X_1, X_2, X_3, X_4).$$

$$\begin{bmatrix} X_1, X_2, X_3 \\ X_2, X_3, X_4 \\ X_1, X_2, X_4 \\ X_2, X_2, X_4 \end{bmatrix} \longrightarrow$$ Repeat until one feature is left or performance decrease is same across all features