

Regression

Supervised learning →

→ We're trying to predict some value given some data.

→ These systems are built using prior labeled data

→ for each X_t , there is an associated Y_t : $\rightarrow \{X_t, Y_t\}_{t=1}^N$

* Setting up data

$\frac{2}{3}$ training $\frac{1}{3}$ testing

- build system using training data,
- test using testing data.
- Compare system against others
- find upper limit of error using this model



training/testing data are from same distribution

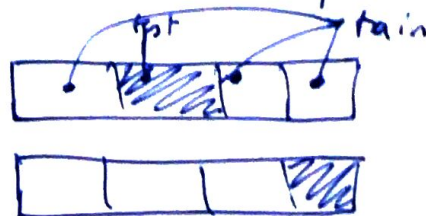
What IF we don't have this much data??

→ we do Cross Validation

we do several training/testing runs while keeping track of errors.
Classifier ↓ Statistics

* S-fold Cross Validation:

- (divide data into S-parts, train on S-4 & test remaining part) S times
- (if sample is really small, we build system on $N-1$ samples & test on just one sample) N times



Learning function:

no noise? $y = f(z)$
data = Z but usually, $X \in Z + \epsilon$
∴ we get,

$g(x) \approx f(z)$
model ideal learnⁿ

Over/UNDER FITTING

→ Our model may give too specific or too broad results initially based on the choice of training data.