FoDA : Cross

L13 : Validation

$\text{Data} \quad (X_{ig}) \quad X \in \mathbb{R}^{n \times d} \quad (d=1) \quad g \in \mathbb{R}^n$

$$\hat{y_i} = M_\alpha(x_i) \longrightarrow y_i$$

poly nominial
model

residual $\hat{y} - y_i$

$$M_\alpha^{(p)}(x) = \sum_{g=0}^{p} \alpha_j x^p = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \dots$$

$$SSE = \sum_i (r_i)^2$$



degree 1 fit

new data

degree 2 fit

degree 3 fit

bias=0

min
variance
sol n

degree 4 fit

degree 5 fit

degree 8 fit

overfitting

Goal: Make good predictions
on new unseen data.

modeling as residual $r_i = \left| M_\alpha(x_i) - g_i \right|$

"generalization"

---

- Collect new data
  $\hookrightarrow$ try on this new data

- "Save" some data for testing

Data

assumption each $(x_i, y_i)$ iid $f$

$X$   $x_i$   $y_i$   $Y$

$X_{train}$   $Y_{train}$

$\alpha = (X_{tr}^T X_{tr})^{-1} X_{tr}^T y_{tr}$

$\downarrow$

then eval

$X_{test}$   $y_{test}$

random hold out

$SSE((X_{test}, y_{test}), \alpha_{train})$

$= \sum_{(x_i, y_i) \in (X_{test}, y_{test})} (M_{\alpha_{train}}(x_i) - y_i)^2$

How well will work on new data? $\sqrt{\frac{1}{|X_{test}|} SSE((X_{test}, y_{test}), \alpha_{train})}$

RMSE

# How large should the test set be?

## Common test size

- 10%
- 33%

---

Evaluate expected value of error by averaging $n_{test}$ observations.

↳ CLT

⟶ unbiased

⟶ variance $\quad \dfrac{Var}{n_{test}}$

more data
less test
percentage

---

more complex model
more test size

# What is cross-validation used for?

- See how model _generalizes_ to new data.

or

- to select a _parameter_ in model (ex. $p$ in $M_\alpha^{(p)}$)

$$p^* = \arg\min_{p \in [1 \ldots 8]} SSE\left((x_{test}, y_{test}), M_{\alpha_{train}}^{(p)}\right)$$

choose best
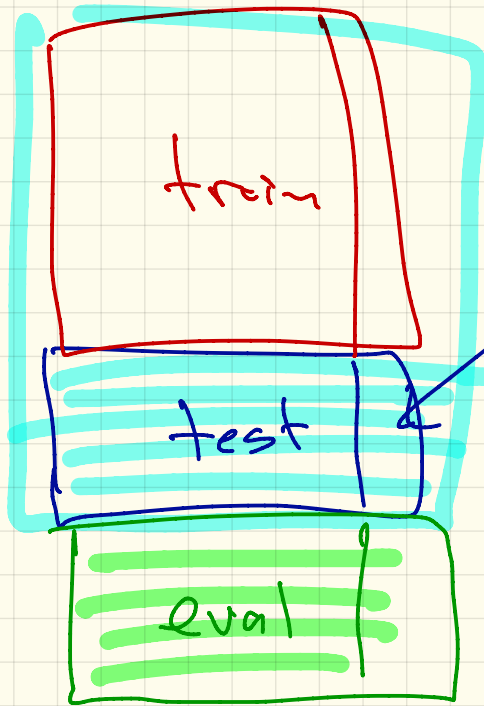
train 8 model

not both

train

dev set

test

← choose
param

eval

← eval
generalizab

OR?

choose   param
          +
evaluate generalizuo

If your data is small,
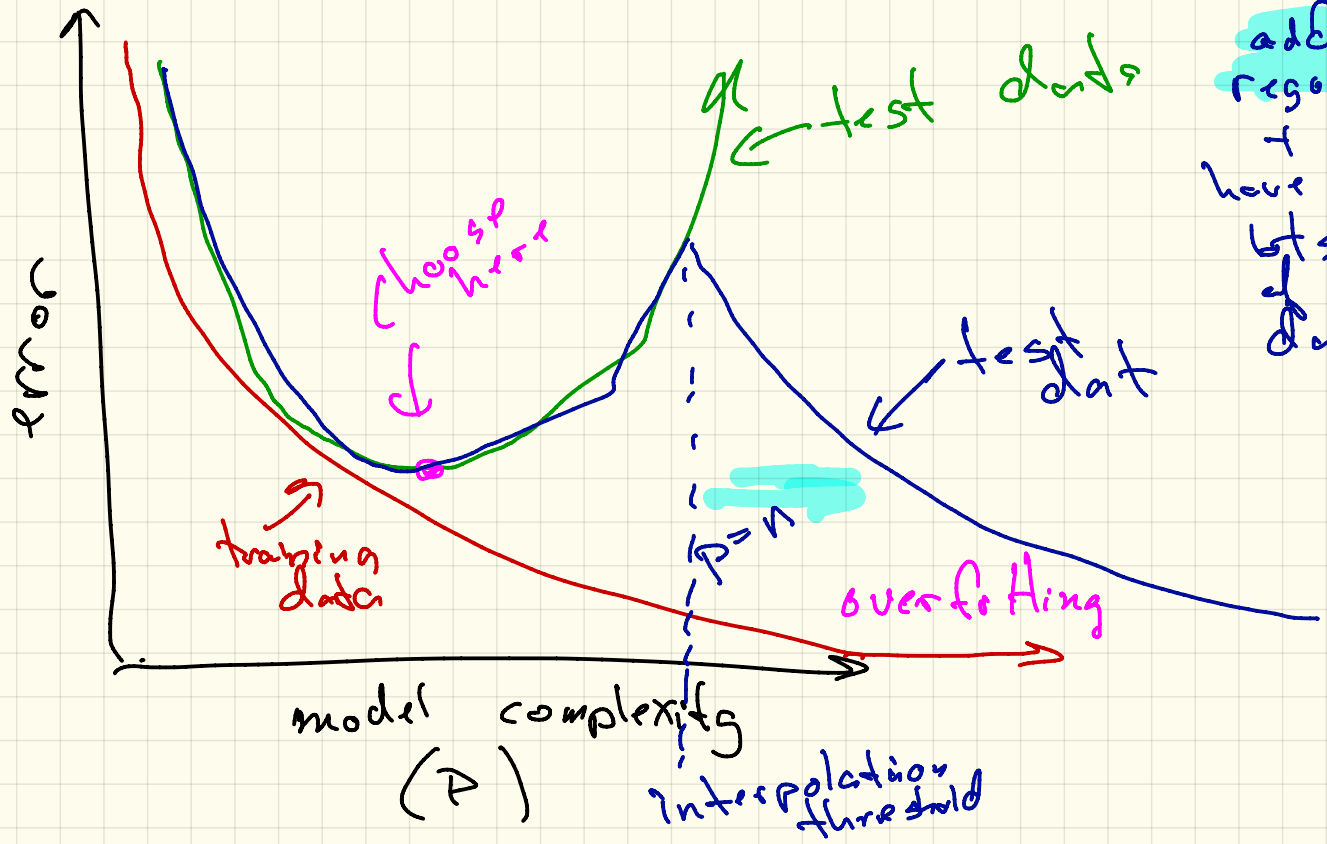don't want to waste data on test.

Leave-one-out CV.

test set size 1
but try all

# Double Descent



error

model complexity

(P)

test data

choose here

test data

training data

$p = n$

overfitting

Interpolation threshold

if you add regularization + have lots! of data