

# EngSci 721

Inverse Problems and Learning From Data

Oliver Maclaren ([oliver.maclaren@auckland.ac.nz](mailto:oliver.maclaren@auckland.ac.nz))

## 1. Basic concepts [5 lectures + 1 Tutorial]

Forward vs inverse problems. Well-posed vs ill-posed problems. Algebra and calculus of inverse problems (left and right inverses, generalised and pseudo inverses, resolution operators, matrix calculus). Representing higher dimensional problems (image data etc).

## 2. Instability and regularisation in linear and nonlinear problems [6 lectures + 1 Tutorial]

Instability and related issues for generalised inverses. Introduction to regularisation and trade-offs. Tikhonov regularisation. Higher-order Tikhonov regularisation. Sparsity and regularisation using different norms. Truncated singular value decomposition. Iterative regularisation, including stochastic/mini-batch gradient descent.

## 3. Further topics [3 lectures + 1 Tutorial]

Regularisation parameter choice, including statistical and machine learning views of regularisation. Confidence sets for linear and nonlinear models. Physics-informed machine learning and neural networks.

# Module overview

Inverse Problems and Learning From Data (*Oliver Maclaren*)

[~14 lectures/3 tutorials]

## Lecture 12: Regularisation parameter choice

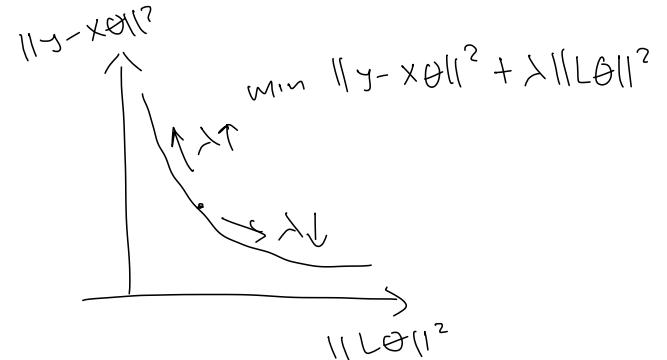
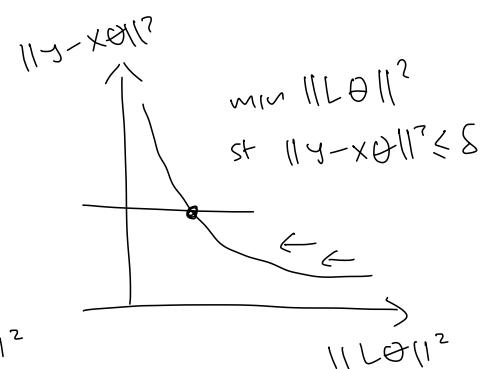
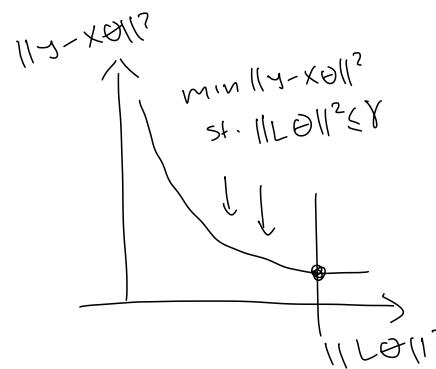
Topics:

- Discrepancy principles
- Cross-validation
- Further reading

## Eng Sci 721 : Lecture 12.

- o Regularisation parameter choice methods
  - Discrepancy approaches
  - Cross-validation
  - Further reading

recall three ways to think about  
trade-off curves:



These are 'equivalent' in the sense  
that they give the same sol<sup>n</sup> for  
appropriate choice of  $\lambda, \delta, \gamma$

## Discrepancy principles

The general idea of a discrepancy principle is to solve eg

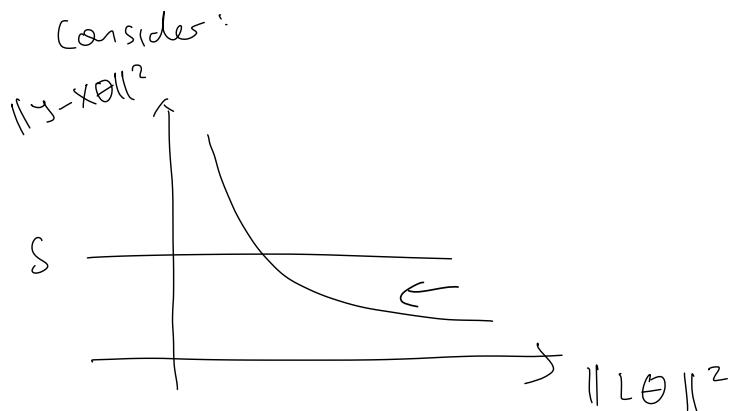
$$\begin{cases} \min \|L\theta\|^2 & \text{model norm term} \\ \text{s.t. } \|y - X\theta\|^2 \leq s & \text{data norm term} \end{cases}$$

for some choice of s

Assuming the constraint is binding  
this means we enforce  $\|y - X\theta\|^2 = s$   
for some s.

Different methods of choosing s exist. Here is one that has a nice interpretation (imo!)

↪ 'simplest adequate/acceptable fitting model'



In some cases we may have info on s, eg if we have a statistical model  $y_i = x_i \theta^{\text{true}} + e_i$  (capitals: RVs)

where the distribution of  $e_i$  is known, eg

$$e_i \sim N(0, \sigma^2), E = [e_1, e_2, \dots]^T$$

→ although we don't know  $e_i$ , we hence know the distribution of  $\|E\|^2$  from the distribution of  $e_i$  ↗

## Chi-squared distribution

Give a series of standard normal random variables, i.e.  $E_i \sim N(0, 1)$   
 $\uparrow$  'standard'

The distribution of

$$Q = \sum_{i=1}^k E_i^2 \quad \left\{ \begin{array}{l} \text{ie sum of squared} \\ \text{errors} \end{array} \right.$$

is called the chi-square / chi-squared distribution with  $k$  'degrees of freedom'

Write  $Q \sim \chi_k^2$   $\uparrow$  number of errors.

If we assume  $y_i = X_i \theta^{\text{true}} + \varepsilon_i$  for  $i = 1, \dots, m$

where  $\varepsilon_i \sim N(0, \sigma^2)$ , then

$$\frac{E_i = Y_i - X_i \theta^{\text{true}}}{\sigma} \quad \left| \begin{array}{l} \text{have dist. } N(0, 1) \\ \text{so} \end{array} \right.$$

$$\Rightarrow \frac{\frac{1}{\sigma^2} \|Y - X\theta^{\text{true}}\|^2}{m} \sim \chi_m \quad \left| \begin{array}{l} \text{standard deviation} \\ \text{random vars here} \end{array} \right.$$

## Hypothesis test inversion & confidence intervals

We don't know  $\theta^{\text{true}}$  but we can test

$$H_0: \theta^{\text{true}} = \theta$$

for any  $\theta$ !

Logic:  $\nearrow$   
 (Fisher)

Suppose data is/are a 'typical' realisation,  $y^{\text{obs}}$ , such that the realised sum of squared errors lies in the 95% region of  $\frac{1}{\sigma^2} \|Y - X\theta^{\text{true}}\|^2$

- If  $\frac{1}{\sigma^2} \|y^{\text{obs}} - X\theta\|^2$  does not lie in the 95% region of  $\frac{1}{\sigma^2} \|Y - X\theta\|^2$  then  $\theta$  cannot be equal to  $\theta^{\text{true}}$   
 $\rightarrow$  rule it out!  
 $\rightarrow$  otherwise call non-rejected  
 $\hookrightarrow$  not the same as 'accepted' as unique answer: ill-posed!

- Can show set of all non-rejected parameters forms a confidence set for  $\theta^{\text{true}}$

$\hookrightarrow$  is a procedure for constructing random sets that will trap true value 95% of the time

## Simplest acceptable fit models

The previous procedure gives us a set of non-rejected (compatible) models.

↳ 'Lack of uniqueness solved' by accepting a set of models

We might stop there! But can also ask for a single representative

→ In general there will be no most complicated model consistent with the data (unbounded complexity)

→ There often is a 'simplest' (in the sense of our regularisation) model though!

↳ gives a lower bound on complexity (Donoho)

Popper! { ↳ simple models = falsifiable  
                          ↳ (can be ruled out)

Also Vapnik... { ↳ complex models = non-falsifiable,  
                          though may be true!

### Procedure:

- Given  $m$  data points
- Find 95% upper limit for  $\hat{x}_m$  distribution (max misfit)
- Find min complexity (regularisation) model fitting no worse than above

Can also e.g. use 50% (median) as a more 'central' point estimate

↳ more typically used approach.

↳ Convenient fact: for large  $m$ , the median of  $\hat{x}_m$  is approx  $\underline{m}$ .

→ choose  $\left[ \|\mathbf{y}^{\text{obs}} - \mathbf{X}\theta\|^2 \approx \sigma^2 \cdot m \right]$  for 'median' discrep. model.

## Discrepancy approaches

### Pros

- statistically well-justified } as interval
- simple interpretation } method

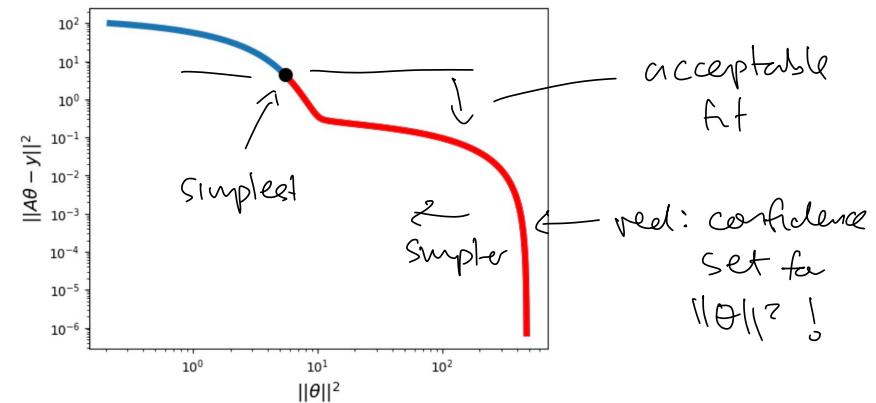
### Cons

- need to know  $\sigma$  or noise model in general
- tends to oversmooth when used as a point estimate  
(remember it's a lower bound!)

Even when taking median tho.

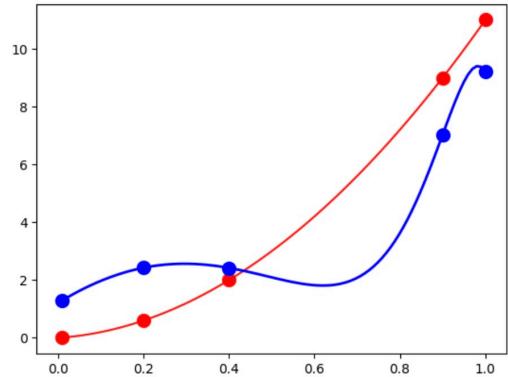
Example: polynomial ( $x^4$ ) regression, known noise level ( $\sigma=1$ )

- Tikhonov  $\|A\theta\|^2$
- loop over each  $\lambda$ , calc norms
- Then:

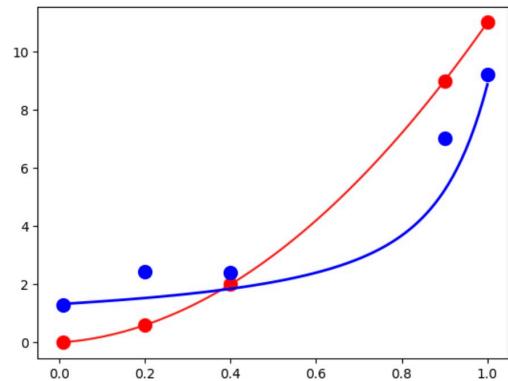


Note: even median quite oversmoothed as point estimate

## Example cont'd



pseudo  
inverse  
sol  
(blue)



discrepancy  
method  
(over  
smoothed  
but treat  
as  
lower  
bound)

## Cross-validation & prediction

Cross-validation & its many variants are also widely used to choose regularisation parameters in inverse problems, statistics & machine learning

As before we aim to 'correct' the naive approach of just considering the models that just fit the given — here, 'training' — data

Why? Firstly,  $\|y^{\text{obs}} - X\theta^{\text{fit}}\|^2$ , where  $\theta^{\text{fit}}$  is the model

from minimising  $\|y^{\text{obs}} - X\theta\|^2$  over  $\theta$ , is known / can be shown to be a biased estimate of

$$\boxed{\text{Err} = \mathbb{E}\left[\|Y - X\theta^{\text{true}}\|^2\right]}$$

(can easily generalise to nonlinear)

which is the expected prediction error,

where  $\theta^{\text{true}}$  is fixed,  $Y$  is randomly drawn from the distribution independent of  $\theta^{\text{true}}$ ,  $X$  is sometimes random, typically fixed, is the expected error for repeatedly predicting new samples at given  $X$  values using the true model.

## Prediction cont'd

Why care about prediction?

1. Intrinsically reasonable thing to care about!
2. It can be shown that  $\underline{\theta}^{\text{true}}$  minimises the expected error for samples drawn from the distribution associated with  $\underline{\theta}^{\text{true}}$  - as it should be! ↳

$$\arg \min_{\theta} \text{Err}(\theta) = \mathbb{E}_{Y|X; \theta^{\text{true}}} [(Y - X\theta)^2] \\ = \theta^{\text{true}}$$

→ this is a form of statistical 'consistency' property.

But we don't know  $\theta^{\text{true}}$ ! And we only have finite data available:  $\mathcal{E}$  rather than  $\mathbb{E}$

## Bias & variance

Although  $\|y^{\text{obs}} - X\theta^{\text{fit}}\|^2$  where  $\theta^{\text{fit}}$  is a model fit using  $y^{\text{obs}}$  (perhaps via least sq.) is typically biased, it may have low variance.

→ In statistics,  $\boxed{\text{risk} = \text{bias}^2 + \text{variance}}$

& we may accept some bias.

↙ under squared error loss

However, if the model family is very large, eg so that we can exactly fit any data, the bias is very large

↳ always estimate error = 0!

→ hence we may want to construct an (approximately) unbiased estimate of the prediction error for any model, really any model fitting strategy

→ this motivates cross-validation

[more sophisticated versions may track some bias & variance of pred. error estimate!]

Recipe: leave-one-out cross validation

Given a sample  $(x, y)_{i=1, \dots, m}$

& a rule for fitting a model to

data  $\lambda \theta_{\lambda}^{\text{fit}}((x, y)_{i=1, \dots, m})$  where

$\lambda$  controls the allowed model complexity, e.g.

$$\theta_{\lambda}^{\text{fit}} = \arg \min_{\theta} \|y^{\text{obs}} - X\theta\|^2 + \lambda \|L\theta\|^2$$

We define  $\theta_{\lambda}^{\text{fit}, i \neq k}$  as the model fit on

all data except the  $k^{\text{th}}$  pair  $(x, y)_k$

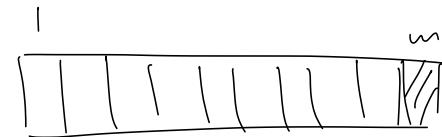
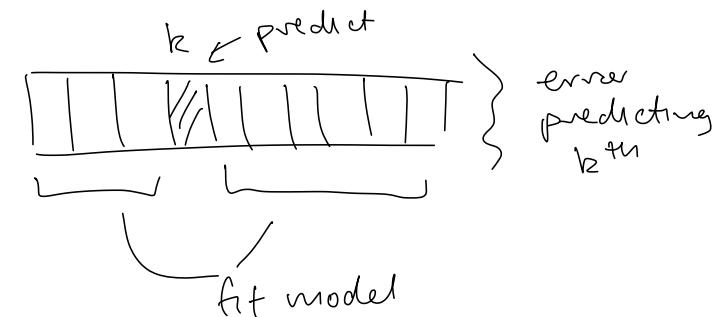
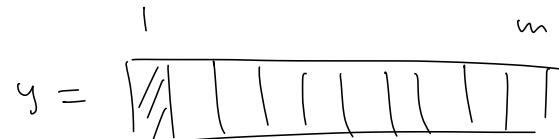
&  $X_{[k, :]}$  as the observed forward map for  
the  $k^{\text{th}}$  observation

& estimate the prediction error of the  
fitting rule by

$$\text{err}(\theta_{\lambda}^{\text{fit}}) = \frac{1}{m} \sum_{k=1}^m (y_k^{\text{obs}} - X_{[k, :]} \theta_{\lambda}^{\text{fit}, i \neq k})$$

i.e. the average error from fitting the  
model on the data with the  $k^{\text{th}}$  data  
point removed & predicting the  $k^{\text{th}}$  data  
point, for  $k=1, \dots, m$  i.e. each data point in  
turn.

Picture:



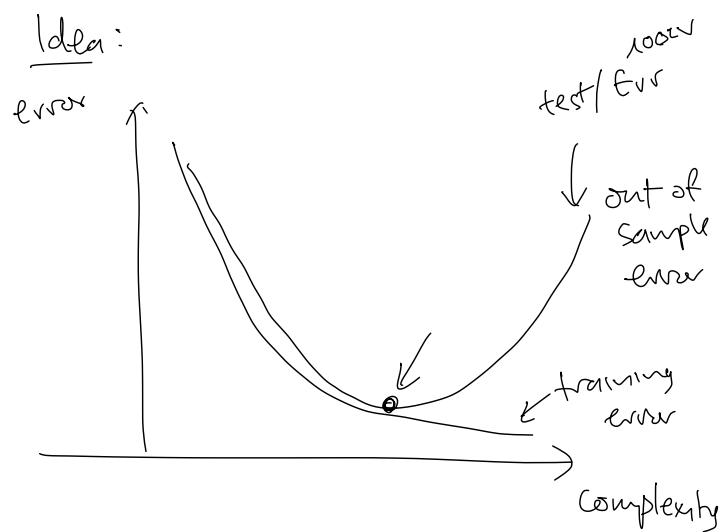
average  
error over  
 $k=1, \dots, m$

## Estimating $\lambda$

Given our estimates for each  $\lambda$ ,  
we choose the  $\lambda$  that minimizes

$$\text{Err}(\hat{\theta}_{\lambda}^{\text{fit}}) = \text{Err}(\lambda)$$

↑  
only depends on  $\lambda$



→ typically refit on whole dataset  
with chosen  $\lambda$  to get final  $\hat{\theta}$  estimate  
(this can cause issues!)

## Notes

- Simple & general, use for any fitting rule
- unbiased estimate of average predictive error for  $m-1$  data points,  $\approx$  average predictive error for  $m$  data points
- averages over both inputs & outputs → error rate for  $\times$  random rather than fixed (ie not error rate for this sample but for rule in general)
- high variance estimate of prediction error!  
↳ can improve with eg  $k$ -fold cross validation or bootstrap versions
- K-fold: leave out  $\approx m/K$  observations at a time instead of one  
→ lower variance, higher bias (typically...)
- computationally expensive!  $K$  fits per  $\lambda$ !

Generalised cross validation: cheaper (for linear)

For linear models with linear regularisation

, it can be written

$$y^{\text{pred}} = X \hat{\theta}_\lambda^{\text{fit}}, \quad \hat{\theta}_\lambda^{\text{fit}}: \begin{matrix} \text{'regularised'} \\ \text{'inverse'} \end{matrix}$$

$$= R_\lambda^\top y^{\text{obs}} \quad \begin{matrix} \text{data space} \\ \text{pseudo projection} \end{matrix}$$

can show a reasonable approximation

to  $\text{err}^{\text{loocv}}(\lambda)$  is

$$\frac{1}{m} \frac{\|y^{\text{obs}} - X \hat{\theta}_\lambda^{\text{fit}}\|^2}{(1 - \text{Trace}(X \hat{\theta}_\lambda^{\text{fit}}))^2} = \frac{1}{m} \frac{\|y^{\text{obs}} - X \hat{\theta}_\lambda^{\text{fit}}\|^2}{(1 - \text{Trace}(R_\lambda))}^2$$

↑  
correction factor

$$\text{like eg } \sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \mu)^2$$

for unbiased variance

↑ Taylor series approx. tends to  $\approx \text{AIC}$

See Aster et al. Inverse Problems

Hastie et al. The elements of statistical learning.

Example: CV attempt

```

thetas = np.zeros((len(alphav), len(theta_true)))
theta_norms = np.zeros((len(alphav), len(y_pobs)))
data_norms = np.zeros((len(alphav), len(y_pobs)))
L = np.eye(len(theta_true))
indices = np.arange(0, len(y_pobs), dtype=int)

for i, alphai in enumerate(alphav):
    for k in range(0, len(y_pobs)):
        # indices
        keep_indices = indices[indices != k]

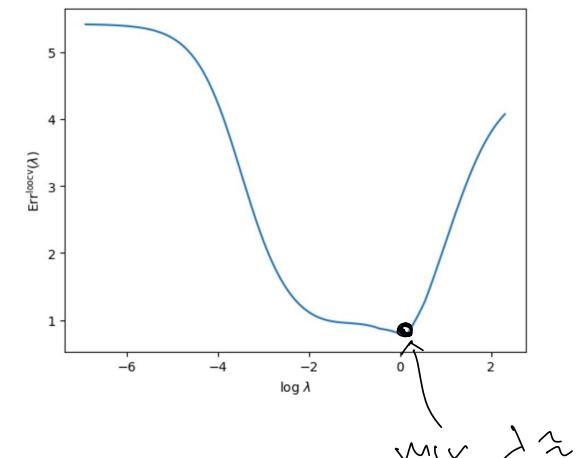
        #augment
        A_poly_aug_k = np.vstack([A_poly_obs[keep_indices, :], alphai * L])
        y_p_aug_k = np.hstack([y_pobs[keep_indices], np.zeros(L.shape[0])])

        #invert
        A_poly_aug_k_pinv = np.linalg.pinv(A_poly_aug_k)
        theta_k = np.dot(A_poly_aug_k_pinv, y_p_aug_k)

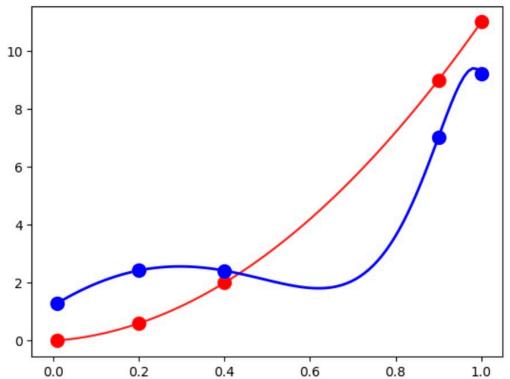
        #calc norms: note prediction on k!
        theta_norms[i, k] = np.linalg.norm(L @ theta_k, 2)
        data_norms[i, k] = np.linalg.norm(y_pobs[k] - np.dot(A_poly_obs[k, :], theta_k))

# average data errors over k for each lambda
data_norms_cv = np.mean(data_norms, axis=1)

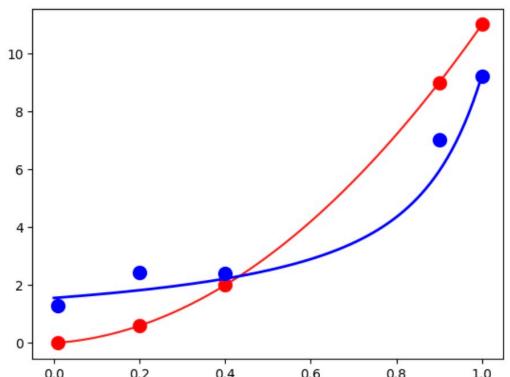
```



## Example cont'd



no regularization  
(min norm  
A<sup>+</sup>)



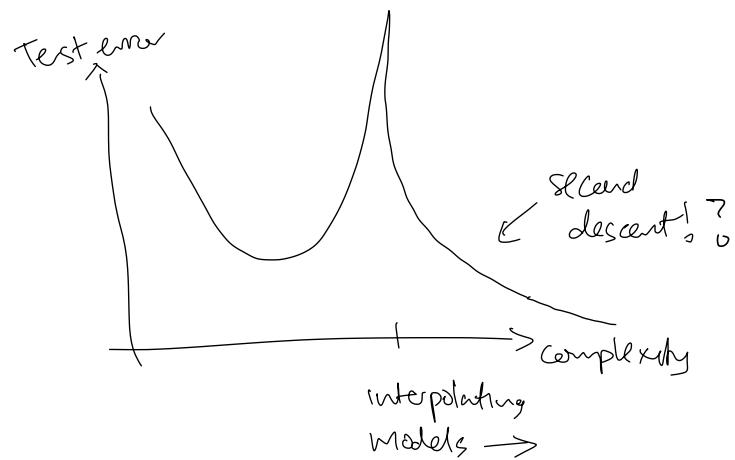
cV +  
Tikhonov

## Further reading

- o Aster et al Inverse Problems
- o Hastie et al The elements of statistical learning
- o Efron (1983) 'Estimating the error rate of a prediction rule'
- o Bates et al (2024) 'Cross-validation: what does it estimate & how well does it do it?'

## Appendix : Double descent?

It has been observed that we can sometimes get plots like



However, imo they usually result  
from plotting vs the wrong measure  
of complexity!

} does demonstrate  
that interpolating  
models can  
generalise  
though!

→ see Assignment 2?

→ Curth et al (2024) 'A U-turn on  
double descent: rethinking  
parameter counting in statistical  
(learning)

etc!