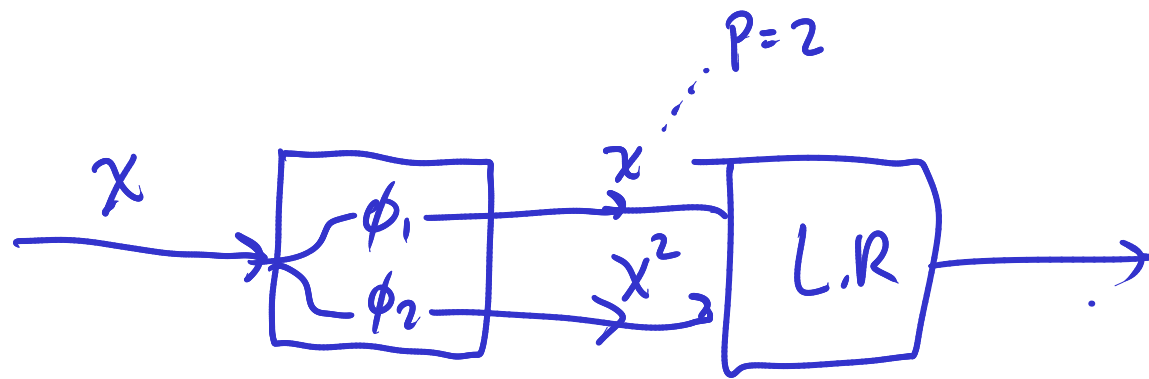
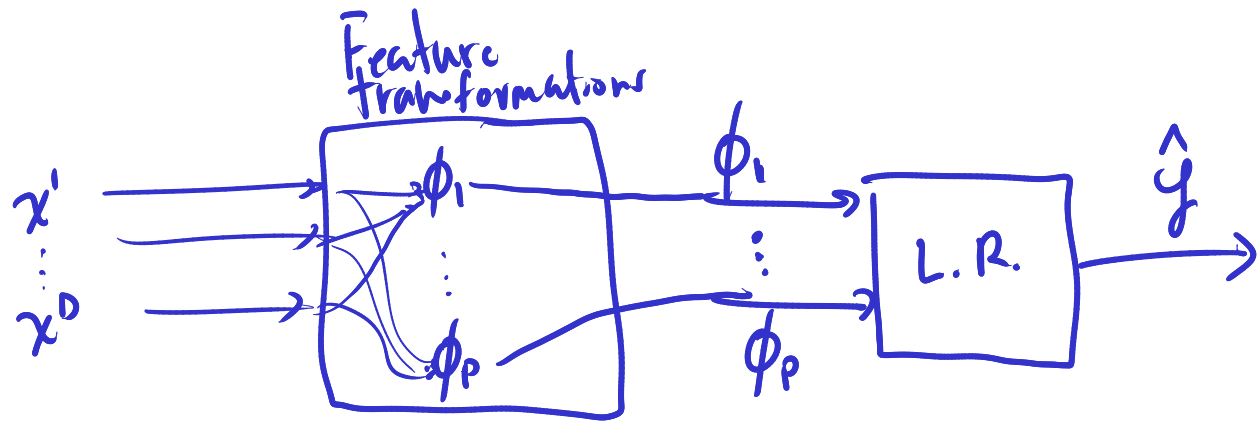




# Statistics and Data Science for Engineers E178 / ME276DS

## Linear regression Part 2

# Recap



## Features

Feature :  $\phi : \mathbb{R}^D \rightarrow \mathbb{R}$ . ... nonlinear transformation.

Takes an input sample and produces a number.

e.g.

$$D=1 : \phi_1(x) = x^2$$

$$D=2 : \phi_1(x^1, x^2) = e^{x^1} \sin x^2$$

$D=1$

## Example: Quadratic fit

$P=2$

trivial ...  $\phi_1(x) = x$   
 $\phi_2(x) = x^2$

$\Rightarrow$

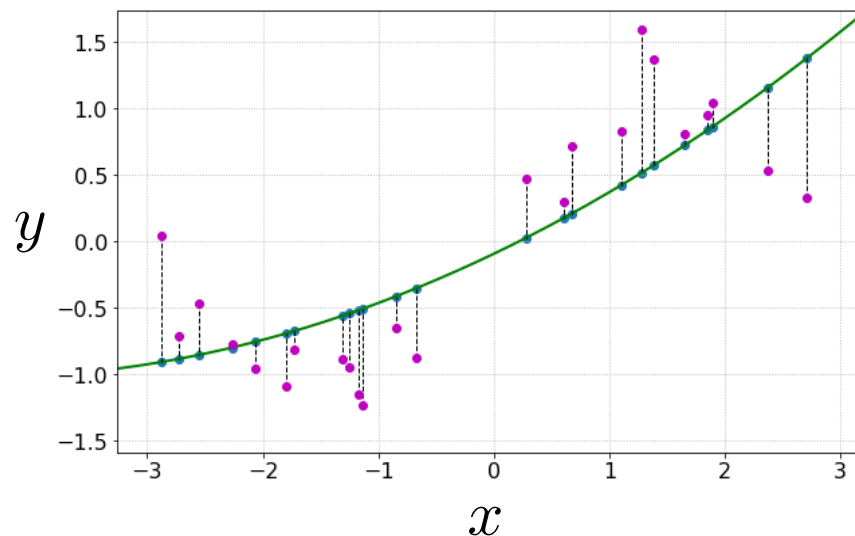
$$h(x_i; \theta_0, \theta_1) = \theta_0 + \phi_1(x_i)\theta_1 + \phi_2(x_i)\theta_2$$
$$= \theta_0 + x_i\theta_1 + x_i^2\theta_2$$

Least squares solution:

$$\hat{\theta}_0 = -0.0935$$

$$\hat{\theta}_1 = 0.4174$$

$$\hat{\theta}_2 = 0.0467$$



In general:  $h(x_i; \theta_0, \theta_1) = \theta_0 + \phi_1(x_i)\theta_1 + \dots + \phi_P(x_i)\theta_P$  ... P features.

Data transformation:  $\Phi = [\phi_1(\mathbf{X}) \quad \dots \quad \phi_P(\mathbf{X})]$

$P > D$ .

$P < D$ .

Solution:

$$\hat{\underline{\theta}}_1 = (\Phi_c^T \Phi_c)^{-1} \Phi_c^T \underline{y}_c$$

$$\hat{\underline{\theta}}_0 = \hat{\mu}_y - \hat{\mu}_\Phi \cdot \hat{\underline{\theta}}_1$$

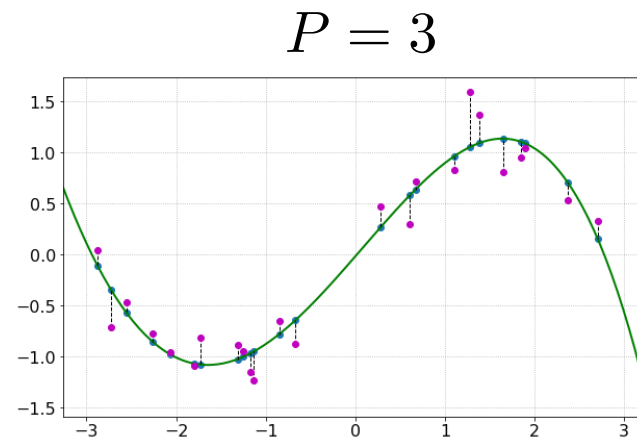
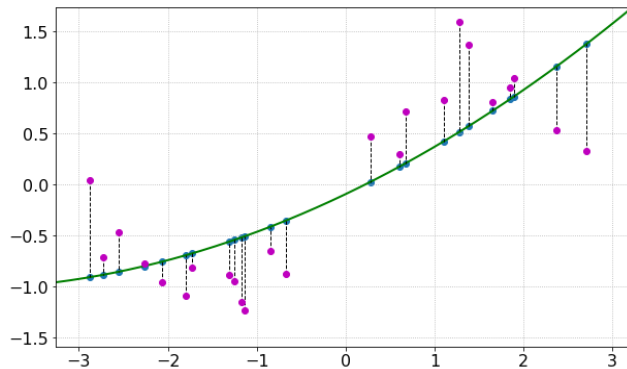
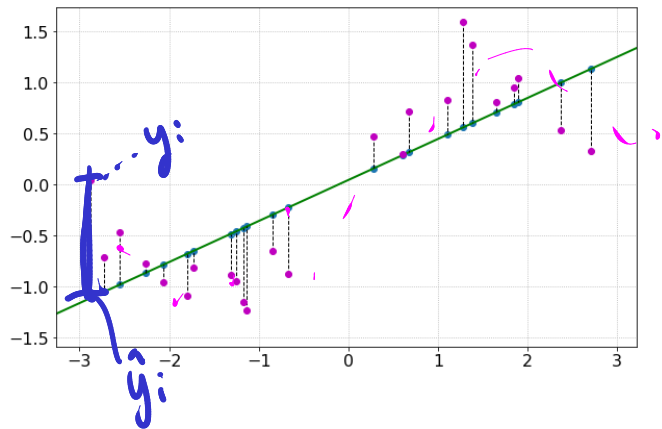
$\phi_1$	$\phi_2$	$\phi_3$	...	$\phi_P$	$y$
0.247746	36.0	266.0	0.247746	88.019654	57.000823
0.179340	37.0	365.0	0.179340	27.211935	22.471227
0.956807	21.0	151.0	0.956807	97.456012	56.357366
0.869653	43.0	437.0	0.869653	43.203221	34.69975
0.825345	47.0	160.0	0.825345	98.933930	64.426163
0.331114	321.0	371.0	0.331114	8.257917	14.218720
0.765523	17.0	364.0	0.765523	96.696783	59.123769
0.956807	21.0	151.0	0.956807	97.456012	52.983470

$\Phi$

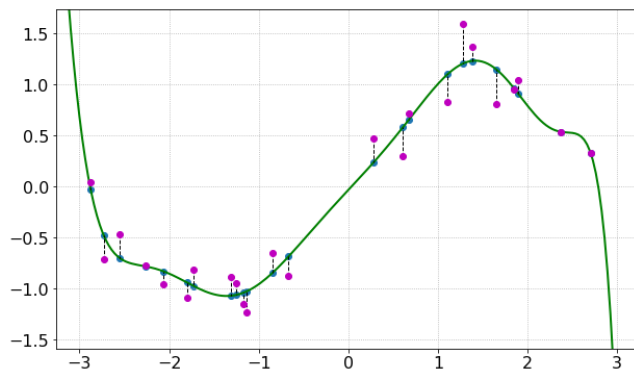
$\mathbf{Y}$

$$\text{minimize } \sum_{P=1} (y_i - \hat{y}_i)^2$$

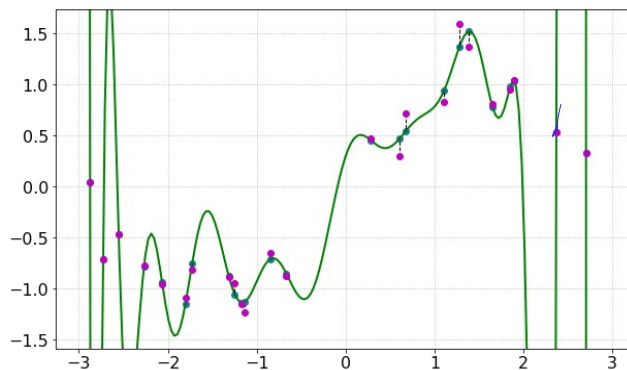
$$\text{MSE} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$



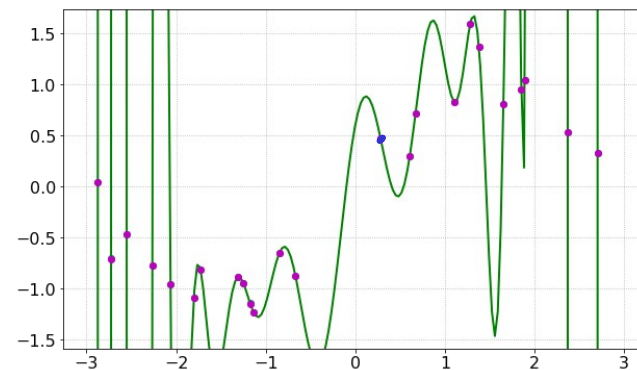
$P=10$



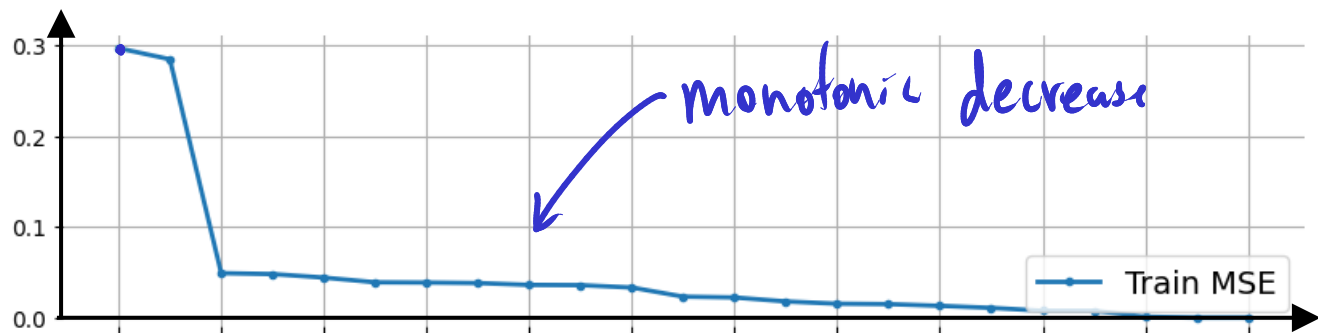
$P=17$



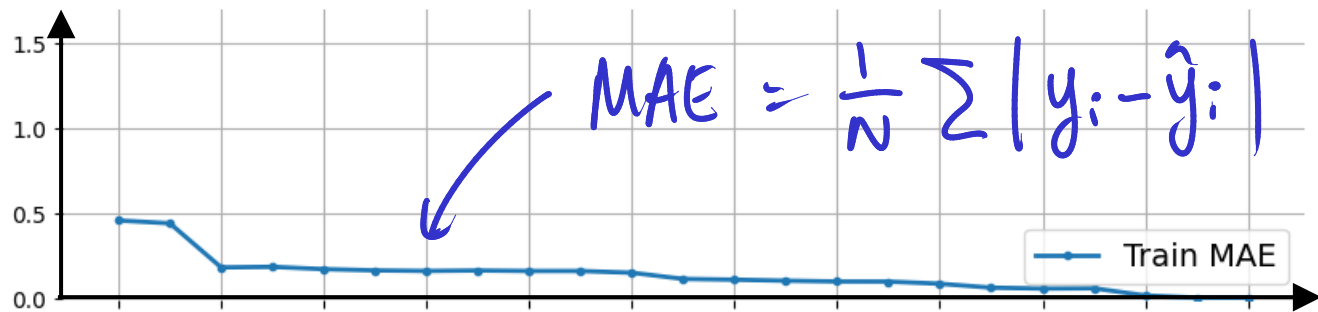
$P=23$



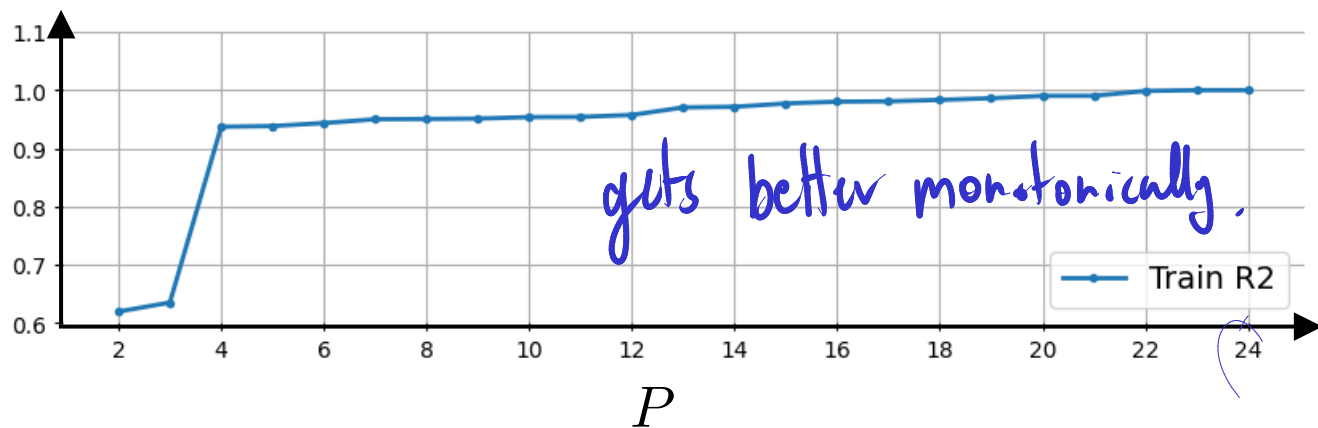
MSE



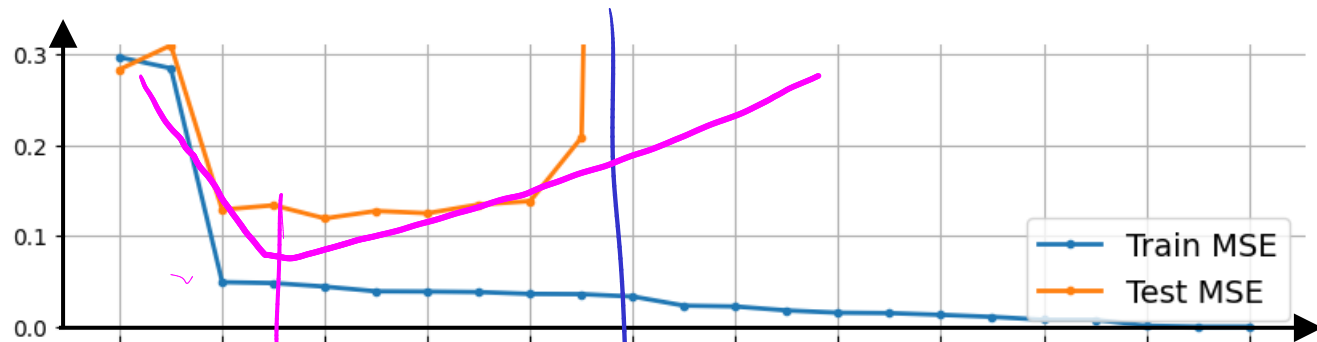
MAE



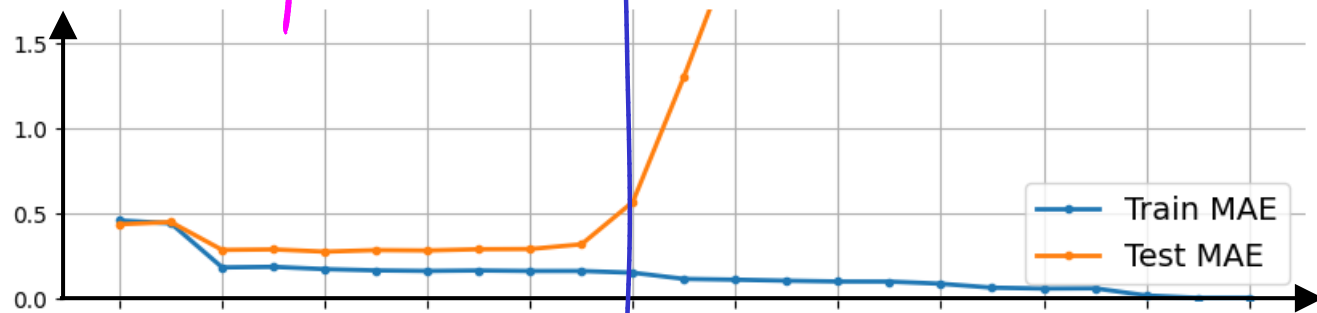
$R^2$



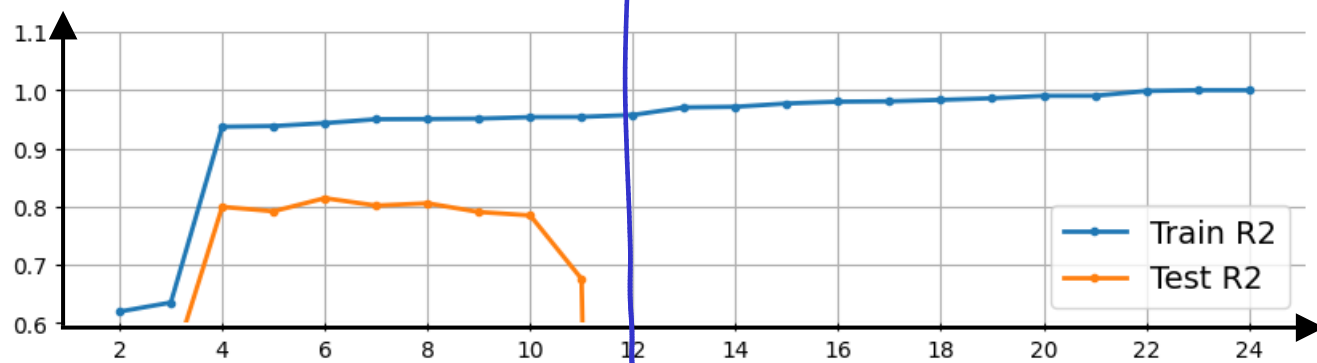
MSE



MAE



$R^2$



$P$

Test data provide a better (unbiased) estimate of the true performance of the model.

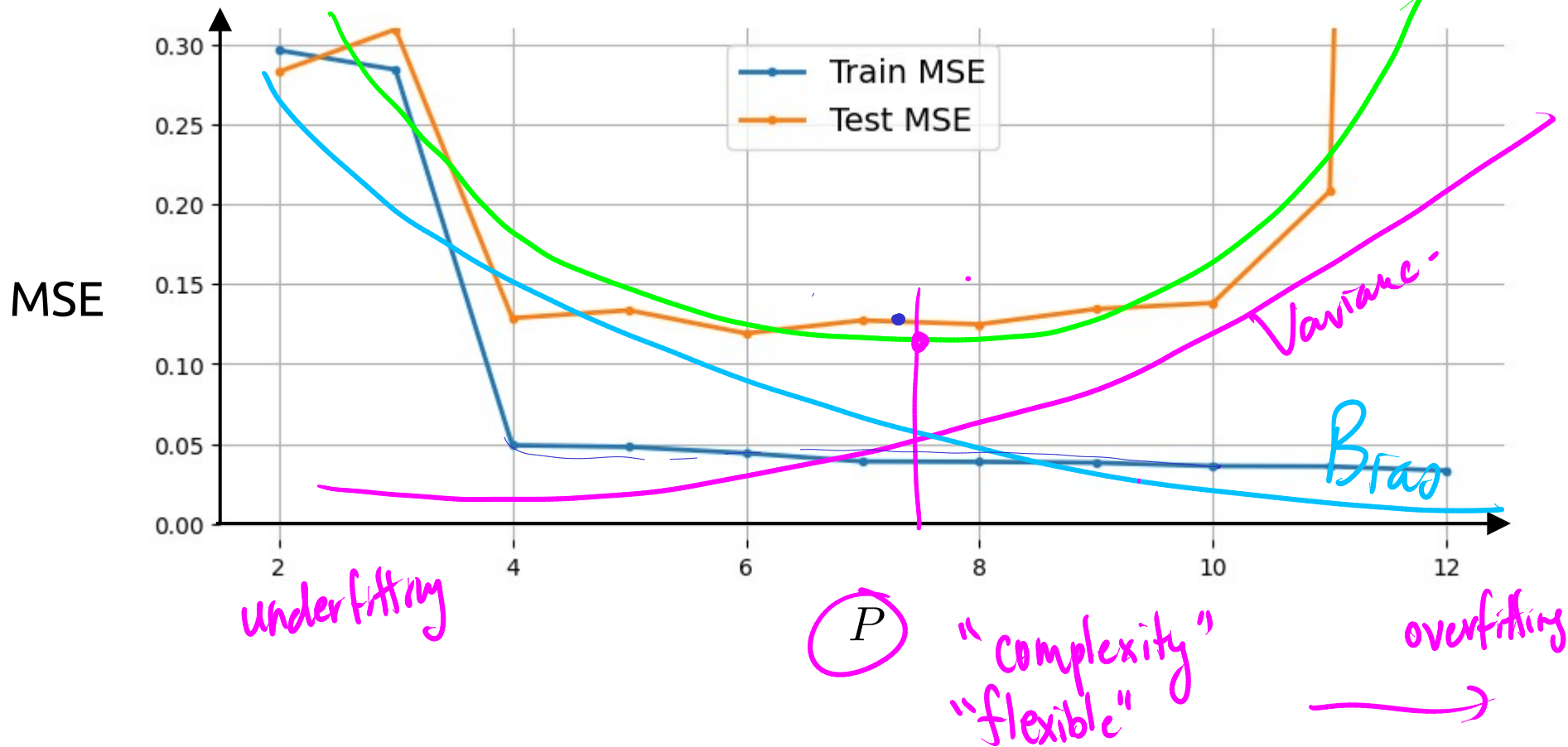
cross-validation.  
also give an unbiased estimate of the true performance.



## Aside: Bias/Variance tradeoff

Bias/Variance decomp.

$$\text{MSE} = \text{Var} + \text{Bias}^2$$



Metrics :

- MSE

- RMSE

- MAE

- MAPE

- $R^2$

... will not detect overfitting.

## Aside: Metrics that penalize complexity

$$\text{AIC} = \frac{1}{\hat{\sigma}^2} \text{MSE}_{\text{train}} + \frac{2}{N} P$$

*penalizes P*

... Akaike information criterion

$$\text{BIC} = \frac{1}{\hat{\sigma}^2} \text{MSE}_{\text{train}} + \frac{\log(N)}{N} P$$

... Bayesian information criterion

where

*C<sub>p</sub> ... Mallows' statistic.*

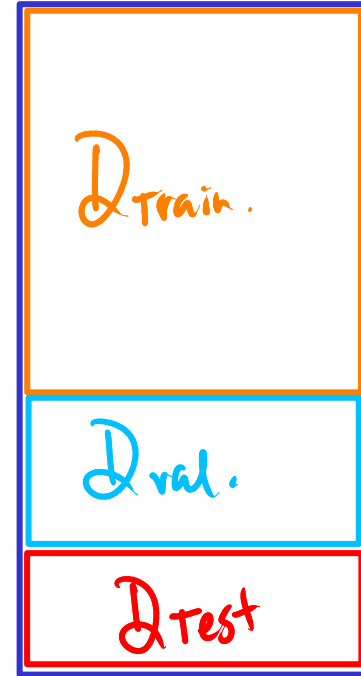
$$\hat{\sigma} = \frac{1}{N - 1 - P} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad \text{for the "full" model}$$

# General feature selection (a.k.a. Model selection)

1. Subset selection.
2. Regularization.

Validation data:

$D =$



# Best subset selection

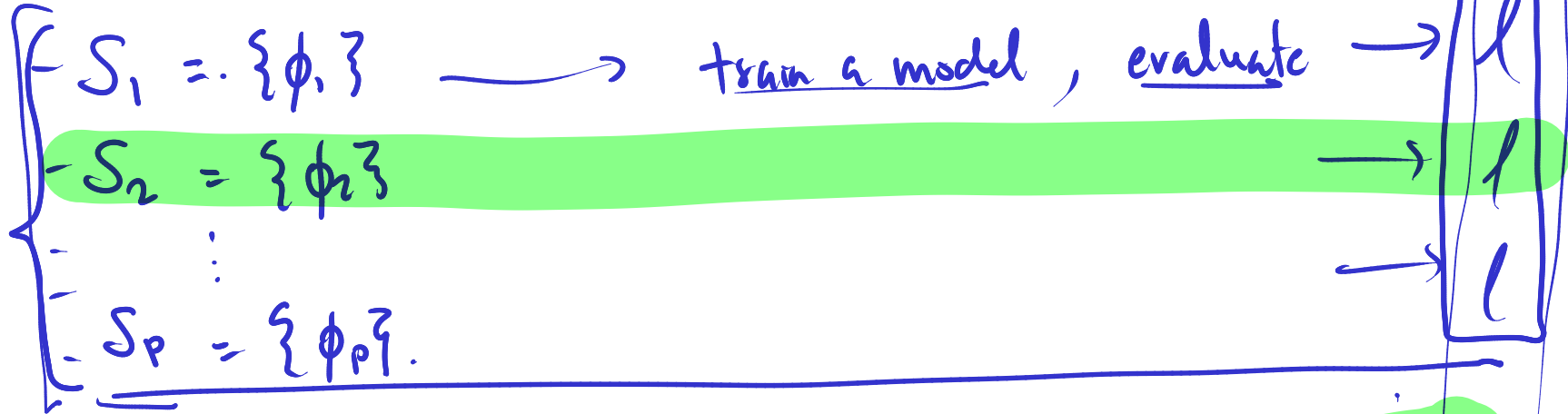
Set of all features

$$\mathcal{P} = \{ \phi_1, \phi_2, \phi_3, \dots, \phi_{P-1}, \phi_P \}$$

$\mathcal{P}$

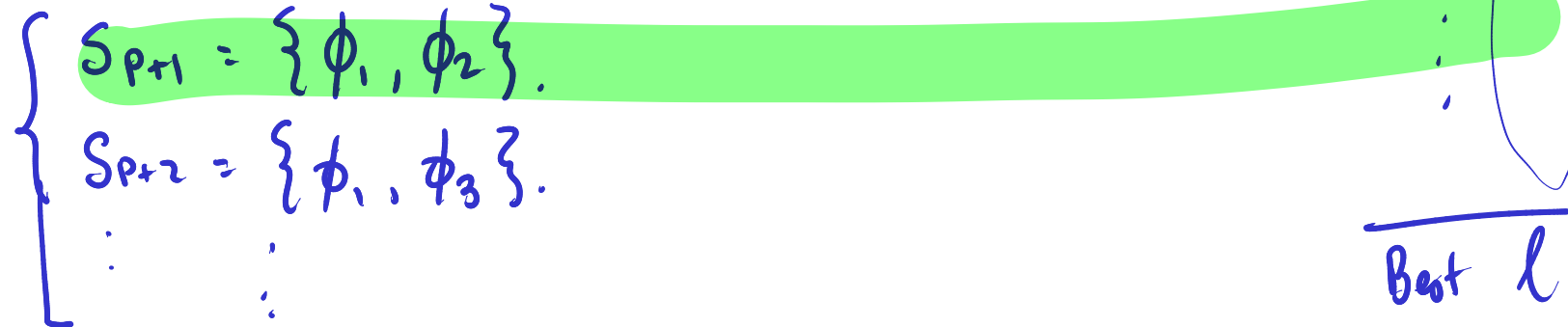
subsets of size 1

$k=1$



subsets size 2

$k=2$



size  $K$

$$\binom{K}{P}$$

$$P=10 \longrightarrow 10^{24}$$

$$P=20 \longrightarrow 10^6$$

$$P=30 \longrightarrow 10^9$$

$K=P$  size  $P$

{

---


$$2^P$$

$P$  things, how many subsets are there?

# Best subset selection

for  $k = 1 \dots P$ :

for  $\mathcal{A}_\kappa$  in {all  $k$ -sized subsets of  $\mathcal{P}$ }:

$$\hat{\theta}_\kappa = \text{train}(\mathcal{A}_\kappa, \mathcal{D}_{\text{train}})$$

$$\ell_\kappa = \text{perf}(\mathcal{A}_\kappa, \hat{\theta}_\kappa, \mathcal{D}_{\text{val}})$$

$$\kappa^* = \text{argbest}(\{\ell_\kappa\})$$

$$\mathcal{S}_k = \mathcal{A}_{\kappa^*}$$

$$\mathcal{S}^* = \text{best of } \{\mathcal{S}_k\}_P$$

$$\hat{\theta}^* = \text{train}(\mathcal{S}^*, \mathcal{D}_{\text{train}})$$

$$\ell^* = \text{perf}(\mathcal{S}^*, \hat{\theta}^*, \mathcal{D}_{\text{test}})$$

# Forward stepwise selection

Assume  $\mathcal{S}_k \subset \mathcal{S}_{k+1}$

$$\# \text{ evaluations} = \sum_{i=1}^P i = \frac{P(P-1)}{2}$$

$$\mathcal{P} = \{ \underline{\phi_1}, \underline{\phi_2}, \phi_3, \dots, \phi_{P-1}, \underline{\phi_P} \}$$

$$\mathcal{S}_0 = \{ \}$$

$$\therefore \mathcal{S}_1 = \{ \quad \phi_3 \quad \} \quad \text{P evaluations.}$$

$$\mathcal{S}_2 = \{ \quad \phi_3 \quad \phi_5 \quad \} \quad \text{P-1 evaluations.}$$

$\vdots$

$$\mathcal{S}_P = \{ \phi_1, \phi_2, \phi_3, \dots, \phi_{P-1}, \phi_P \}$$



# Forward stepwise selection

$$\mathcal{S}_0 = \{\}$$

for  $k = 1 \dots P$ :

for  $\kappa, \phi_p \in \text{enumerate}(\mathcal{P} \setminus \mathcal{S}_{k-1})$

$$\mathcal{A}_\kappa = \mathcal{S}_{k-1} \cup \phi_p$$

$$\hat{\theta}_\kappa = \text{train}(\mathcal{A}_\kappa, \mathcal{D}_{\text{train}})$$

$$\ell_\kappa = \text{perf}(\mathcal{A}_\kappa, \hat{\theta}_\kappa, \mathcal{D}_{\text{val}})$$

$$\kappa^* = \text{argbest}(\{\ell_\kappa\})$$

$$\mathcal{S}_k = \mathcal{A}_{\kappa^*}$$

$$\mathcal{S}^* = \text{best of } \{\mathcal{S}_k\}_P$$

$$\hat{\theta}^* = \text{train}(\mathcal{S}^*, \mathcal{D}_{\text{train}})$$

$$\ell^* = \text{perf}(\mathcal{S}^*, \hat{\theta}^*, \mathcal{D}_{\text{test}})$$

set minus.  
all unused features.  
in stage  $k-1$

```

curlyS = [set() for i in range(P+1)]
ellk = np.full(P+1, np.inf)

for k in range(1, P+1):

    curlyA = [set() for i in range(P-k+1)]
    ellkappa = np.full(P-k+1, np.inf)

    for kappa, phip in enumerate(curlyS[k-1] - curlyS[k]):
        curlyA[kappa] = curlyS[k-1].union({phip})
        theta0hat, thetalhat = train(curlyA[kappa], Dtrain)
        ellkappa[kappa] = perf(curlyA[kappa], theta0hat, thetalhat,
                               Dvalidate)

    kappastar = ellkappa.argmin()
    curlyS[k] = curlyA[kappastar]
    ellk[k] = ellkappa[kappastar]

kstar = ellk.argmin()
Sstar = curlyS[kstar]
theta0star, thetalstar = train(Sstar, Dtrain)
ellstar = perf(Sstar, theta0star, thetalstar, Dtest)

# Store the results
f_ellk = ellk
f_ellstar = ellstar
f_kstar = kstar
    
```

# Backward stepwise selection

Assume  $\mathcal{S}_k \subset \mathcal{S}_{k+1}$

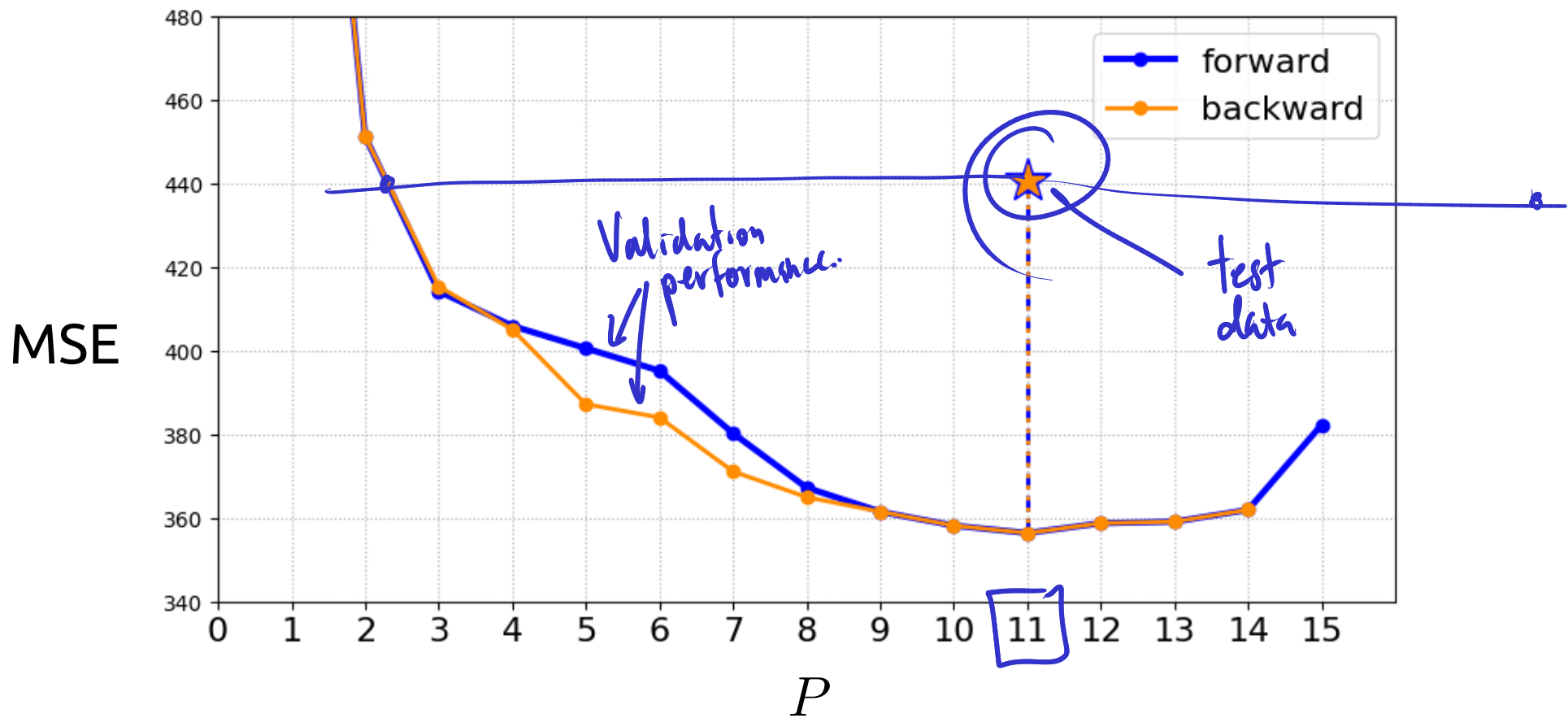
$$\mathcal{S}_P = \{ \phi_1, \phi_2, \phi_3, \dots, \phi_{P-1}, \phi_P \}$$

$$\mathcal{S}_{P-1} = \{ \phi_1, \cancel{\phi_2}, \phi_3, \dots, \phi_{P-1}, \phi_P \} \quad P-1 \text{ evaluations.}$$

$$\mathcal{S}_{P-2} = \{ \phi_1, \cancel{\phi_2}, \phi_3, \dots, \phi_{P-1}, \phi_P \}$$

$\vdots$

$$\mathcal{S}_0 = \{ \}$$



# Parameter shrinkage, a.k.a. Regularization

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left( \underbrace{\sum_{i=1}^N L(y_i, \hat{y}_i)}_{\text{loss minimization}} + \underbrace{\lambda R(\theta)}_{\text{balancing parameter}} \right)$$

penalty on the size of the parameters.  $\theta$ .

$$\lambda \geq 0$$

$\lambda$  small  $\longrightarrow$  overfitting

$\lambda$  large  $\longrightarrow$  less overfitting.

$R$  can be a norm function.  $\begin{cases} \longrightarrow 1 \text{ norm (absolute values) LASSO.} \\ \longrightarrow 2 \text{ norm (Euclidean) Ridge.} \end{cases}$

# Ridge regression

(linear regression

L2.

slope parameters.

$$\hat{\theta}_{\text{ridge}} = \underset{\theta_0 \dots \theta_P}{\operatorname{argmin}} \left( \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^P \theta_j^2 \right)$$

$$= \underset{\theta_0, \underline{\theta}_1}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \hat{\mathbf{Y}} \right\|_2^2 + \lambda \left\| \underline{\theta}_1 \right\|_2^2 \right)$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_P \end{bmatrix}$$

intercept

slope parameters

always invertible  $\lambda > 0$ .

matrix notation.

Solution:

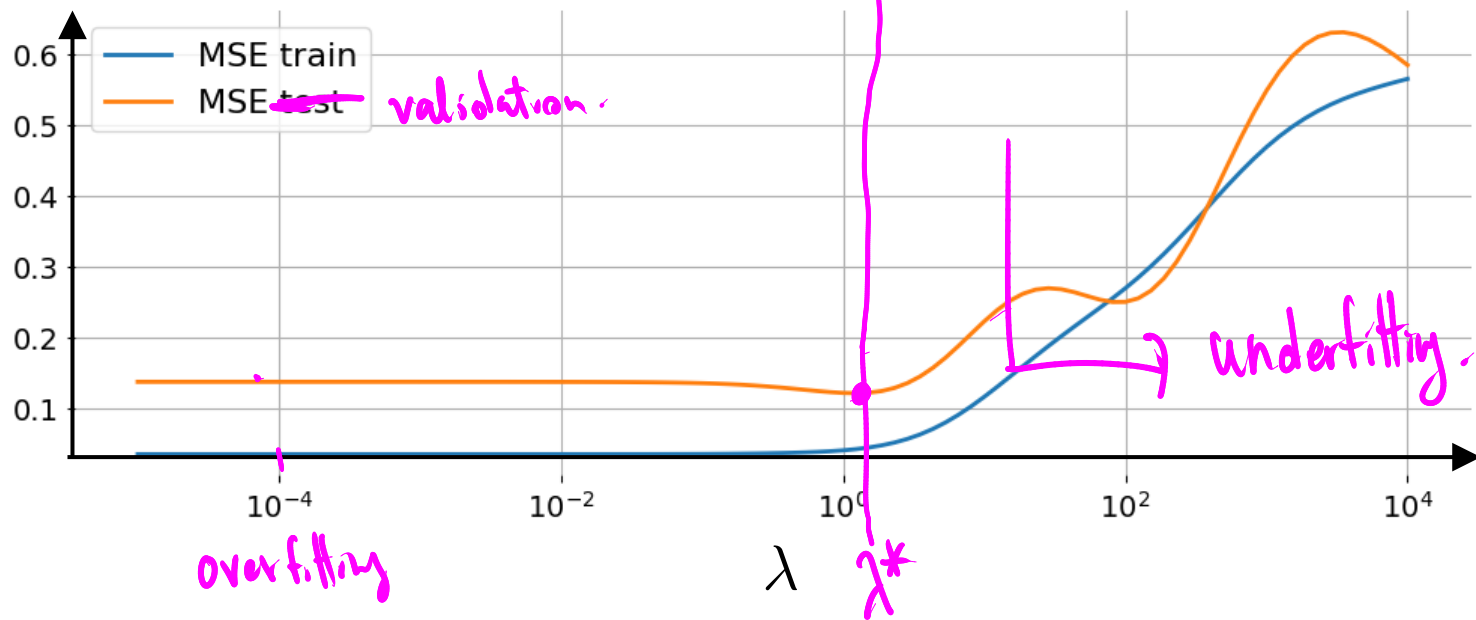
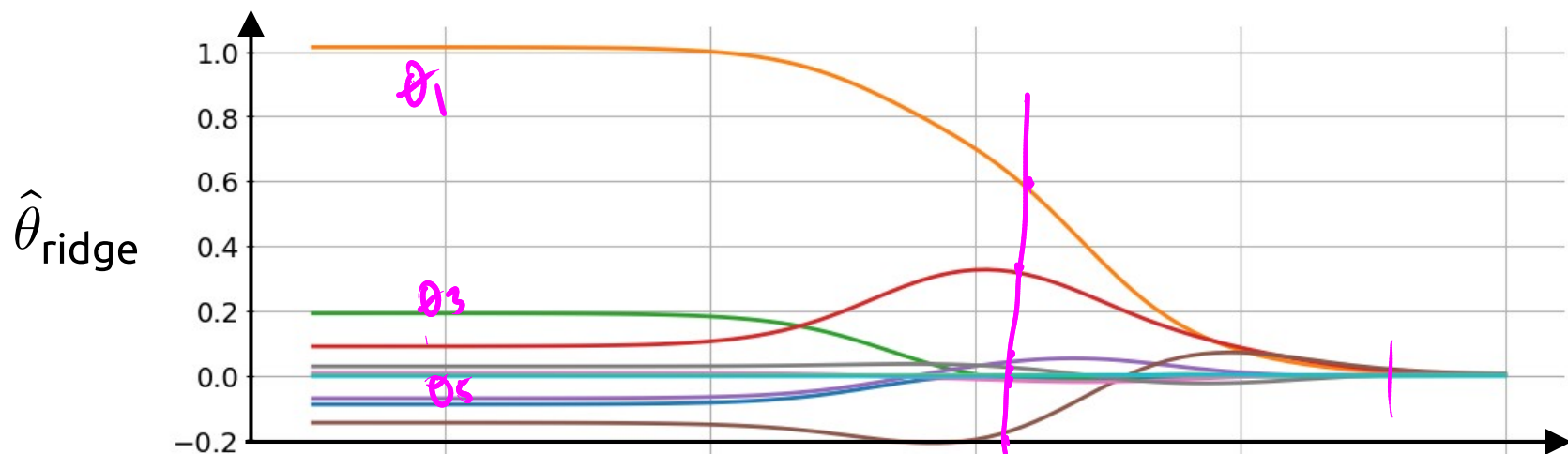
$$\hat{\theta}_0 = \hat{\mu}_Y - \hat{\mu}_X \hat{\theta}_1$$

$$\hat{\theta}_1 = ((\mathbf{X}^c)^T \mathbf{X}^c + \lambda I)^{-1} (\mathbf{X}^c)^T \mathbf{Y}^c$$

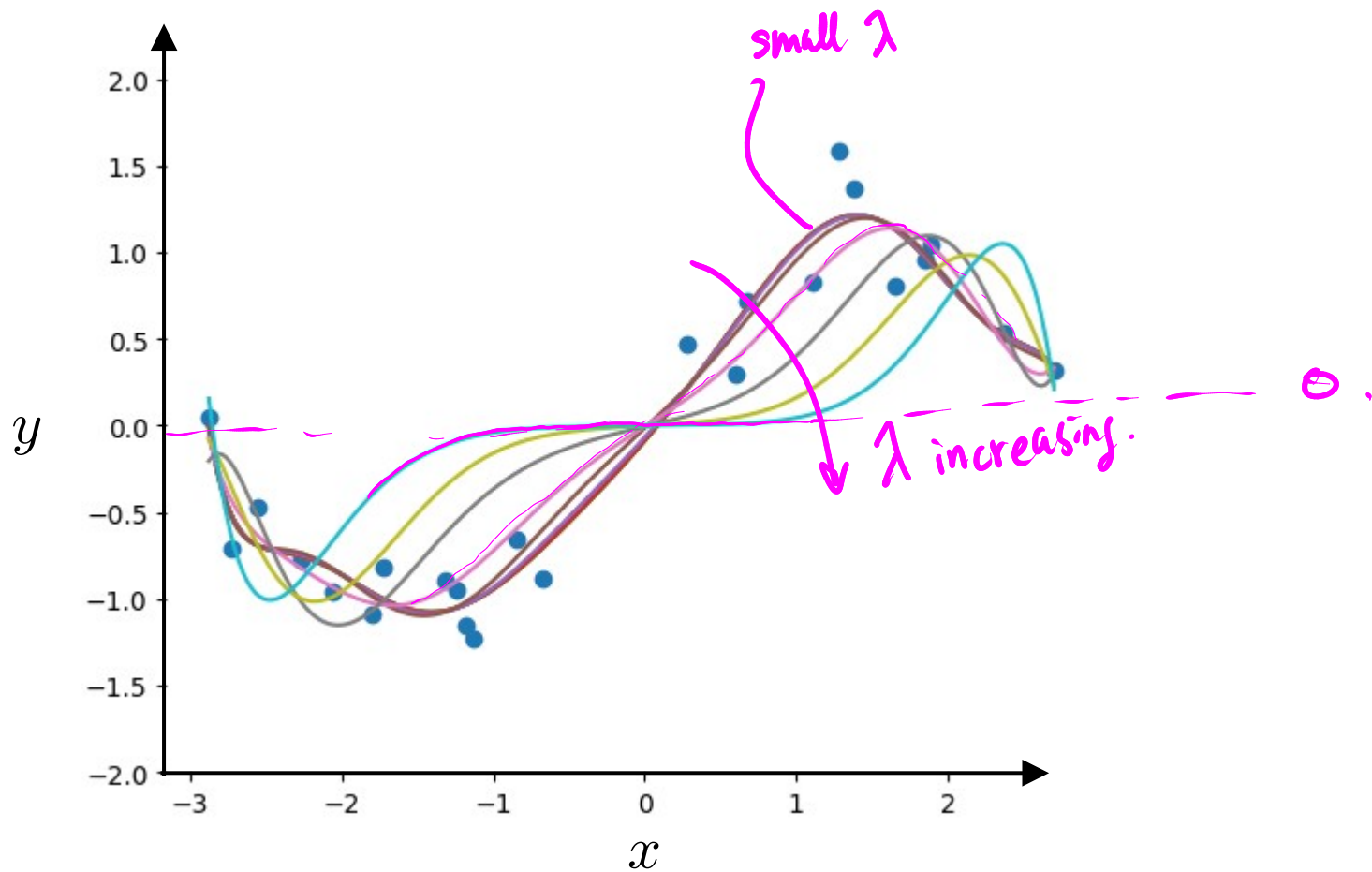
identity matrix.

Before assumed invertible.

$$\hat{\theta}_1 = (\mathbf{X}^{cT} \mathbf{X}^c + \lambda I)^{-1} (\mathbf{X}^c)^T \mathbf{Y}^c$$



# Ridge regularized linear regression



# Lasso regression

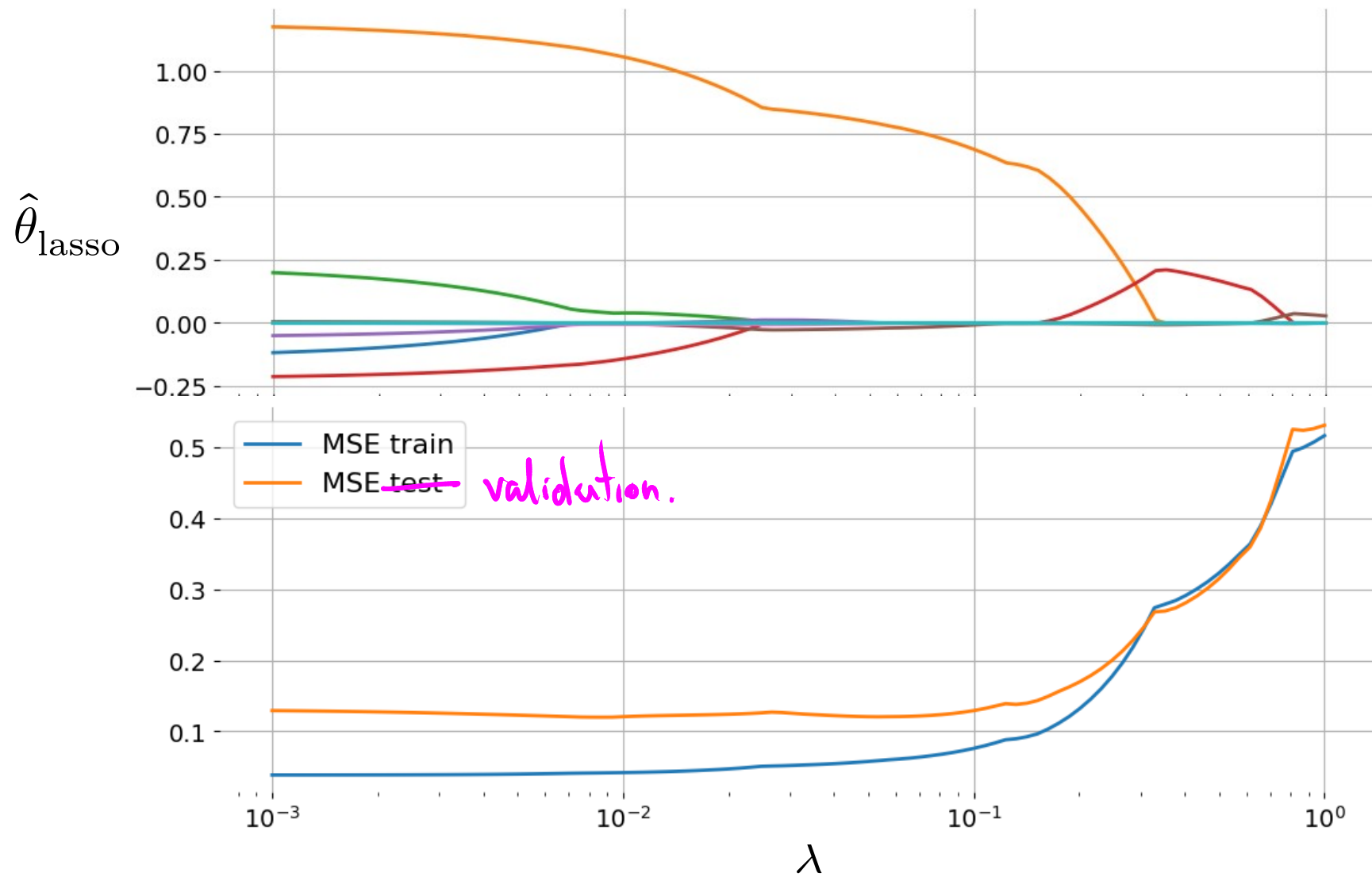
$$\hat{\theta}_{\text{lasso}} = \underset{\theta_0 \dots \theta_P}{\operatorname{argmin}} \left( \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^P |\theta_j| \right)$$

$$= \underset{\theta_0, \underline{\theta}_1}{\operatorname{argmin}} \left( \left\| \mathbf{Y} - \hat{\mathbf{Y}} \right\|_2^2 + \lambda \left\| \underline{\theta}_1 \right\|_1 \right)$$

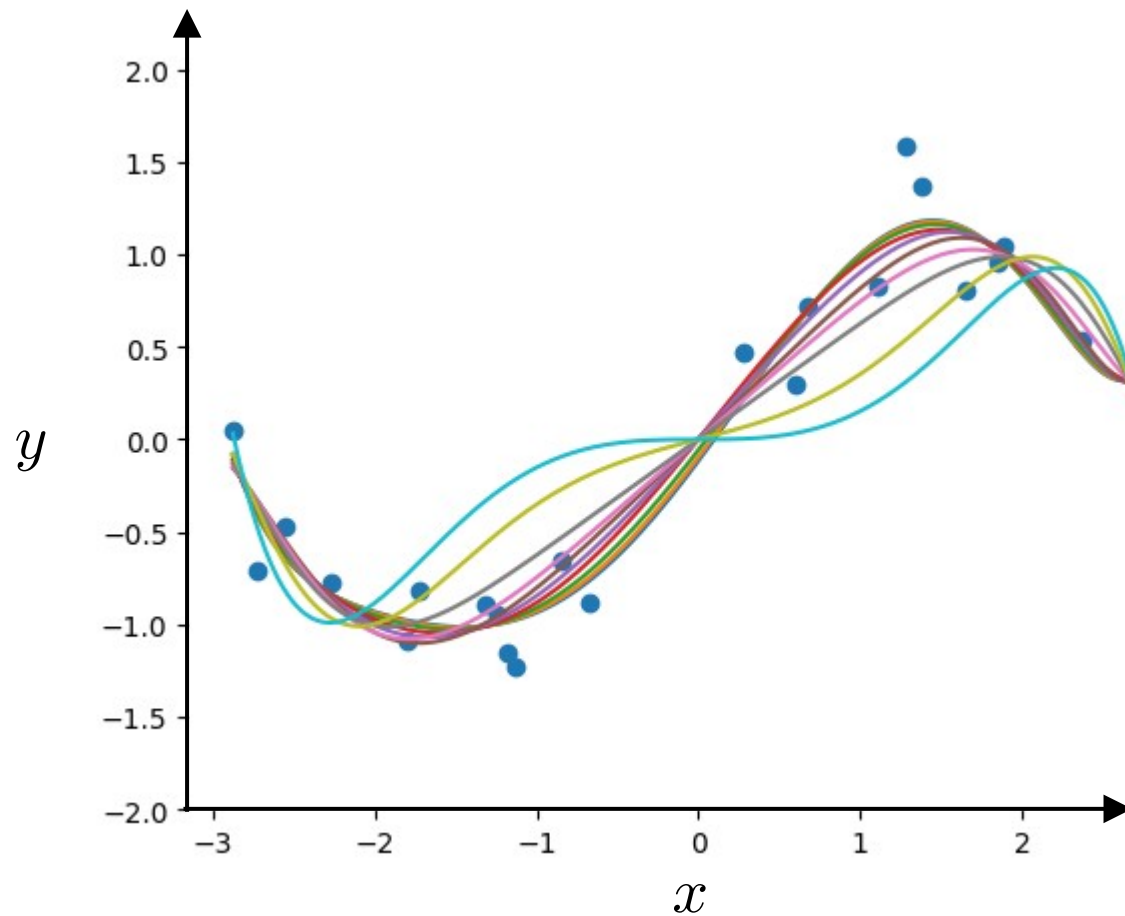
↑ 1 norm.

Creates sparse models. Several  $\theta$ 's = 0.





# Lasso regularized linear regression



## A different perspective on Ridge and Lasso

$$\hat{\theta}_{\text{ridge}} = \underset{\theta_0, \theta_1}{\operatorname{argmin}} \left( \| \mathbf{Y} - \hat{\mathbf{Y}} \|^2 + \lambda \| \theta_1 \|_2^2 \right)$$

$$= \underset{\theta_0, \theta_1}{\operatorname{argmin}} \| \mathbf{Y} - \hat{\mathbf{Y}} \|^2$$

*Lagrange multiplier*

$$\text{s.t. } \| \theta_1 \|_2 \leq t$$

*depend on  $\lambda$*

$$\hat{\theta}_{\text{lasso}} = \underset{\theta_0, \theta_1}{\operatorname{argmin}} \left( \| \mathbf{Y} - \hat{\mathbf{Y}} \|^2 + \lambda \| \theta_1 \|_1 \right)$$

$$= \underset{\theta_0, \theta_1}{\operatorname{argmin}} \| \mathbf{Y} - \hat{\mathbf{Y}} \|^2$$

$$\text{s.t. } \| \theta_1 \|_1 \leq t$$

