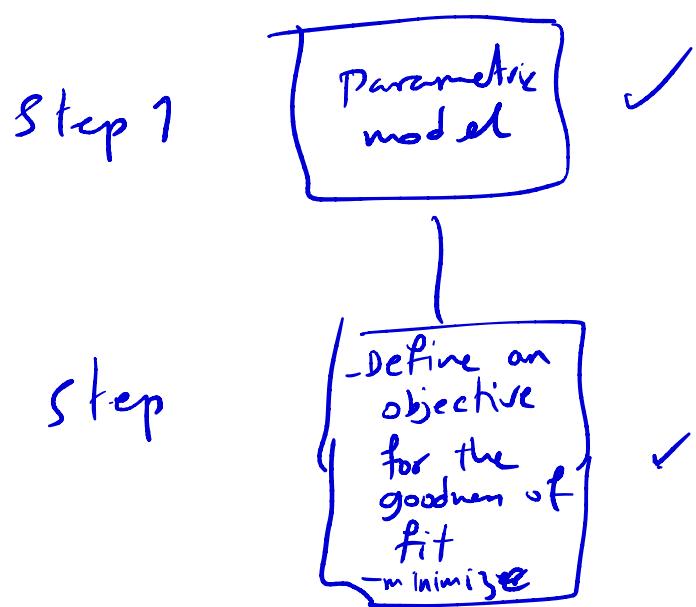


## Lecture 15

P-dimensional linear regression:

$$\text{Data: } \{(x_i, y_i)\}_{i=1}^n \quad x_i \in \mathbb{R}^P$$



Step 1:

$$x_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^P \rightarrow y_i \approx \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \dots + \hat{\beta}_p x_{ip}$$

P+1 parameters  
(to be learned)  
from data

Step: 2:

$$RSS = \sum_{i=1}^n (\text{error}_i)^2$$

$$= \sum_{i=1}^n \left( y_i - \underbrace{\left( \hat{B}_0 + \hat{B}_1 x_{i1} + \dots + \hat{B}_p x_{ip} \right)}_{\substack{\text{true label} \\ \text{prediction}}} \right)^2$$

to find  $\{\hat{B}_j\}_{j=0}^p$  we'll have to  
minimize RSS.

$$RSS = \sum_{i=1}^n (y_i - \vec{x}_i^T \vec{B})^2$$

$$\hat{B} = \begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix} \in \mathbb{R}^{p+1}, \vec{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} \in \mathbb{R}^{p+1}$$

$$RSS = \|\vec{Y} - X \vec{B}\|_2^2$$

$$\boxed{\vec{Y} = \begin{bmatrix} 1, x_{i1}, x_{i2}, \dots, x_{ip} \end{bmatrix}^T, \vec{B} = \begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix}}$$

$\hat{B}_0 + \hat{B}_1 x_{i1} + \dots + \hat{B}_p x_{ip}$   
 $\text{data}_i$

$$\vec{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}_{n \times 1}$$

$$\vec{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}_{n \times (p+1)}$$

$$= \begin{pmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_n^T \end{pmatrix}$$

$$\| \vec{z} \|_2^2 = \sum_{i=1}^n z_i^2$$

$$\vec{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix}$$

$$\underset{\vec{\beta}}{\text{minimize}} \quad RSS(\vec{\beta})$$

$$\underset{\vec{\beta}}{\text{minimize}} \quad \underbrace{\| \vec{Y} - \vec{X} \vec{\beta} \|_2^2}_{RSS(\vec{\beta})}$$

$$\Rightarrow \underbrace{\nabla \text{RSS}(\vec{B})}_{\downarrow} = 0$$

$$-2 \underset{=}{} \vec{X}^T (\vec{y} - \vec{X} \vec{B}) = 0$$

$$\Rightarrow \vec{X}^T \vec{X} \vec{B} = \vec{X}^T \vec{y}$$

$$\Rightarrow \vec{B} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$$

$$\begin{pmatrix} \hat{B}_0 \\ \hat{B}_1 \\ \vdots \\ \hat{B}_p \end{pmatrix}$$

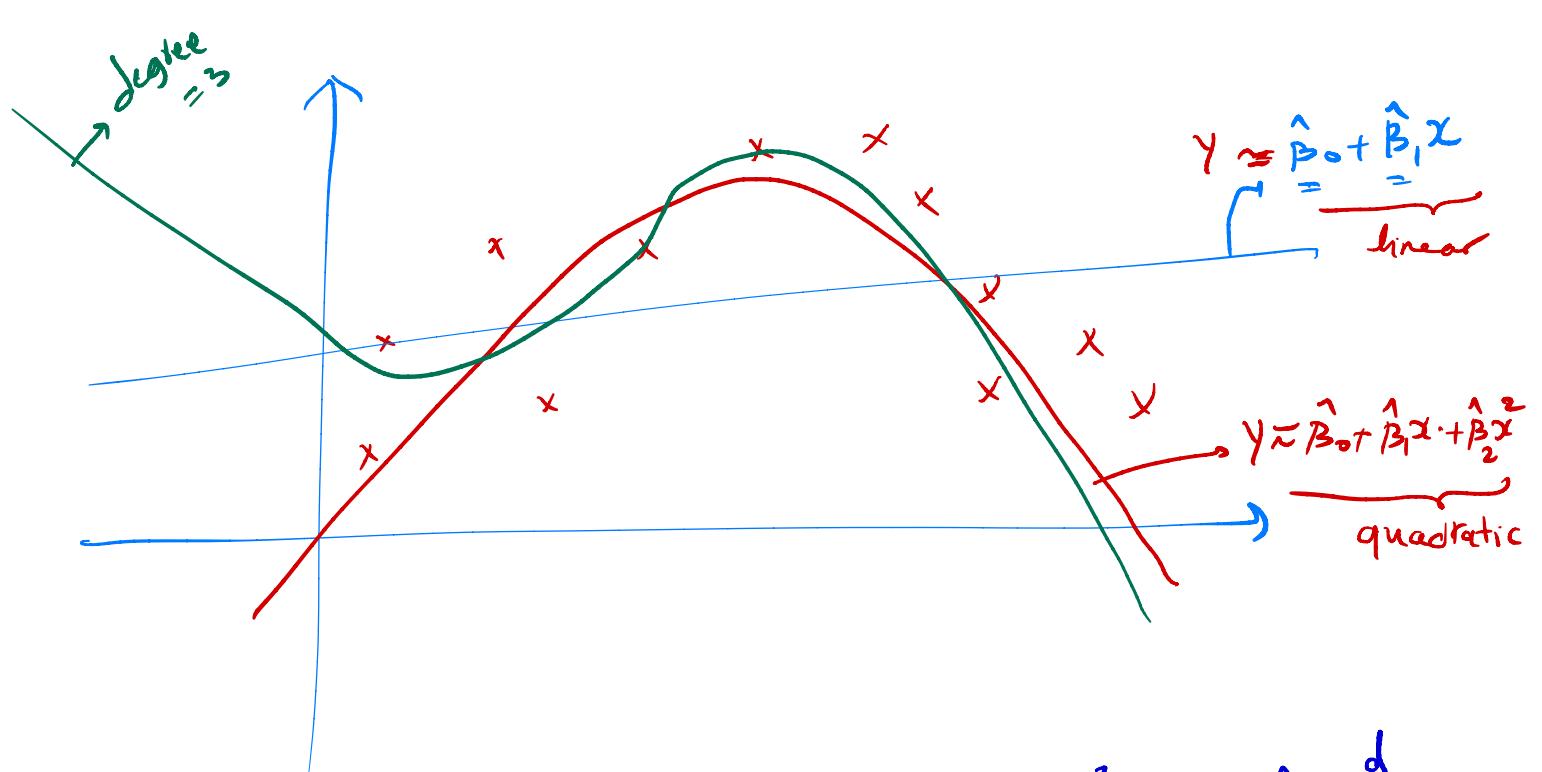
$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{i1} & \cdots & x_{ip} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}$$

i-th

- So far, we've only talked about linear models.  
 Let's now see how we can construct more complex (non-linear) parametric models..

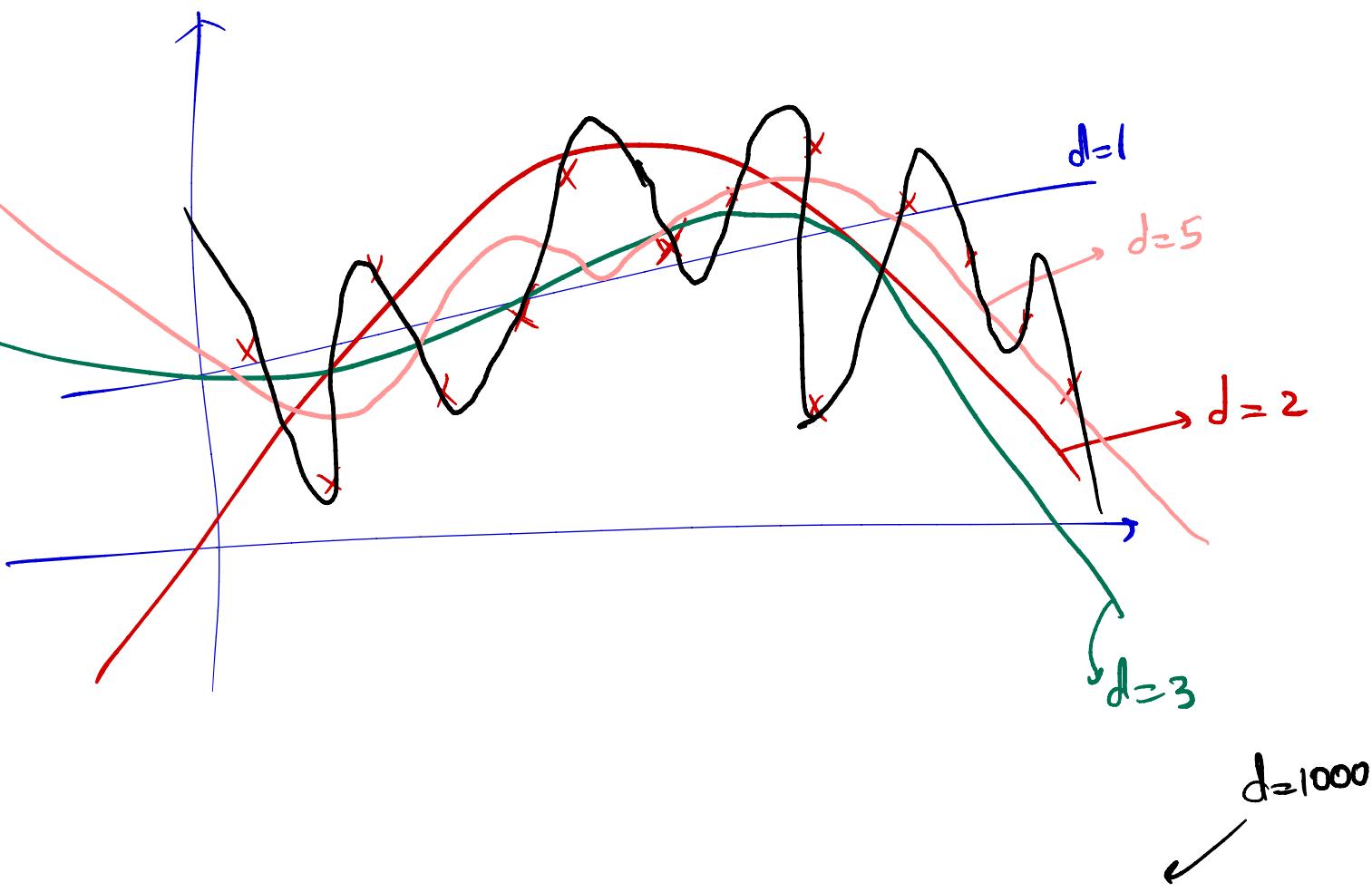
Let's assume  $p=1$ .



polynomial :  $y = p(x) = \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \dots + \hat{\beta}_d x^d$

One way to construct non-linear parametric models is by using polynomial models:

$$\left\{ \begin{array}{l} \text{degree of the model: } d \\ \text{parameters: } \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_d \quad (\text{d+1 parameters}) \end{array} \right.$$



$$y \approx \hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 x^2 + \dots + \hat{\beta}_{1000} x^{1000}$$

Question: Given a polynomial family of degree  $d$ , how do we find the best parameters (the best fit)?

} Data:  $\{(x_i, y_i)\}_{i=1}^n$      $x_i \in \mathbb{R}$ ,  $y_i \in \mathbb{R}$   
 } Step 1: Parametric model: polynomials of degree  $d$ :  $y = \hat{\beta}_0 + \hat{\beta}_1 x + \dots + \hat{\beta}_d x^d$

Step: Define an objective for the goodness of fit:

$$\xrightarrow{\text{objective}} \text{RSS} = \sum_{i=1}^n \text{error}_i^2$$

$$= \sum_{i=1}^n \left( y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\beta}_2 x_i^2 + \dots + \hat{\beta}_d x_i^d) \right)^2$$

this objective is analogous to the multidimensional linear regression:

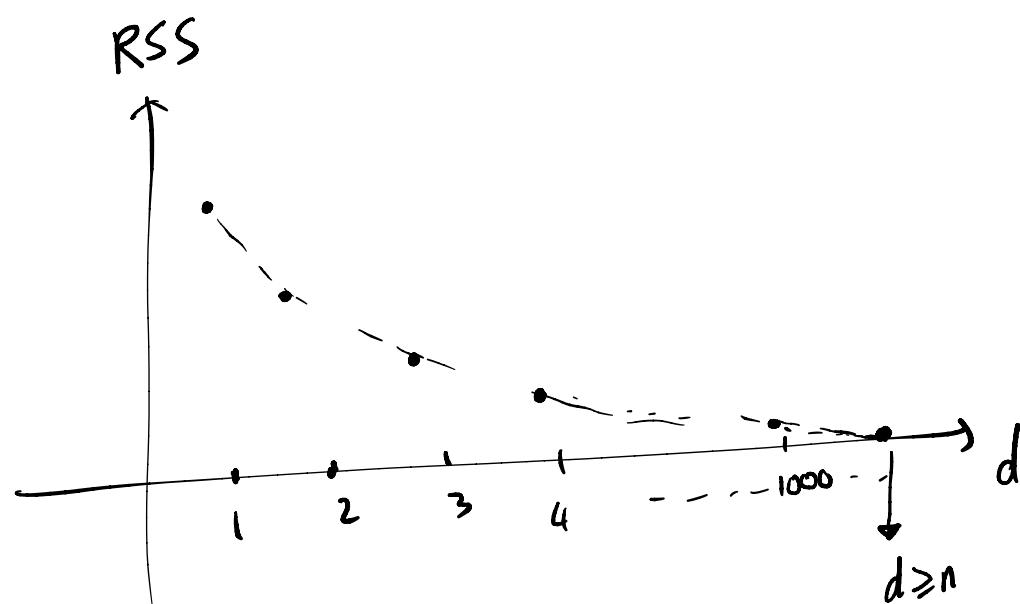
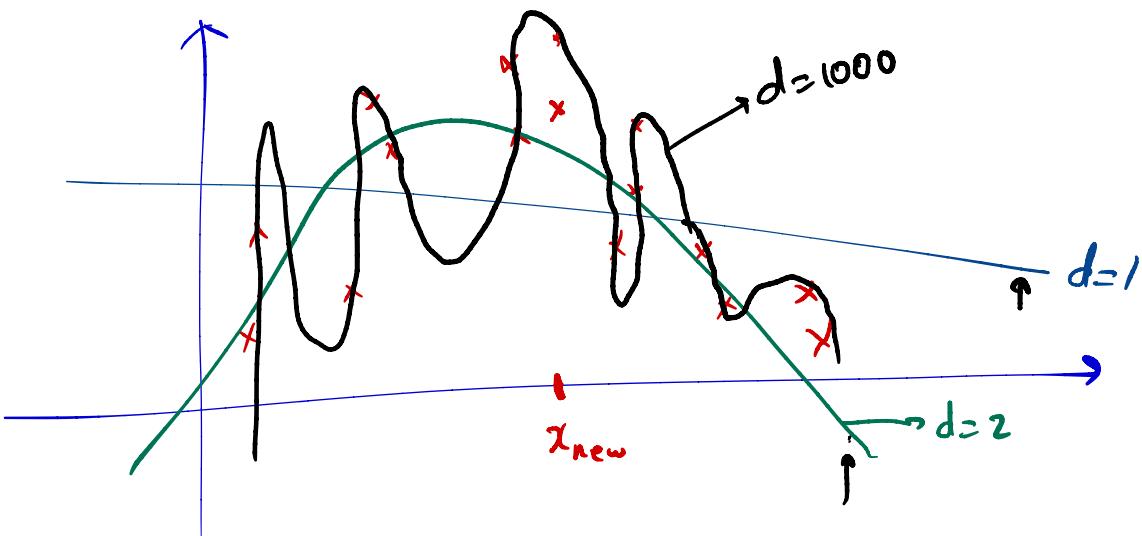
$$\underline{x_i} \in \mathbb{R} \rightarrow \tilde{x}_i = \begin{pmatrix} 1 \\ x_i \\ x_i^2 \\ \vdots \\ x_i^d \end{pmatrix}, \quad \vec{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \vdots \\ \hat{\beta}_d \end{pmatrix}$$

new data point #i

$$\text{RSS} = \sum_{i=1}^n (y_i - \tilde{x}_i^T \vec{\beta})^2$$

$$\Rightarrow \vec{\beta} = (\vec{x}^T \vec{x})^{-1} (\vec{x}^T \vec{y})$$

$$\vec{x} = \begin{pmatrix} 1 & x_1 & x_1^2 & \cdots & x_1^d \\ 1 & x_2 & x_2^2 & \cdots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^d \end{pmatrix}_{n \times (d+1)} \quad \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$



- Remember that the main goal behind learning predictive models was to be able to predict well on new and unseen data points.

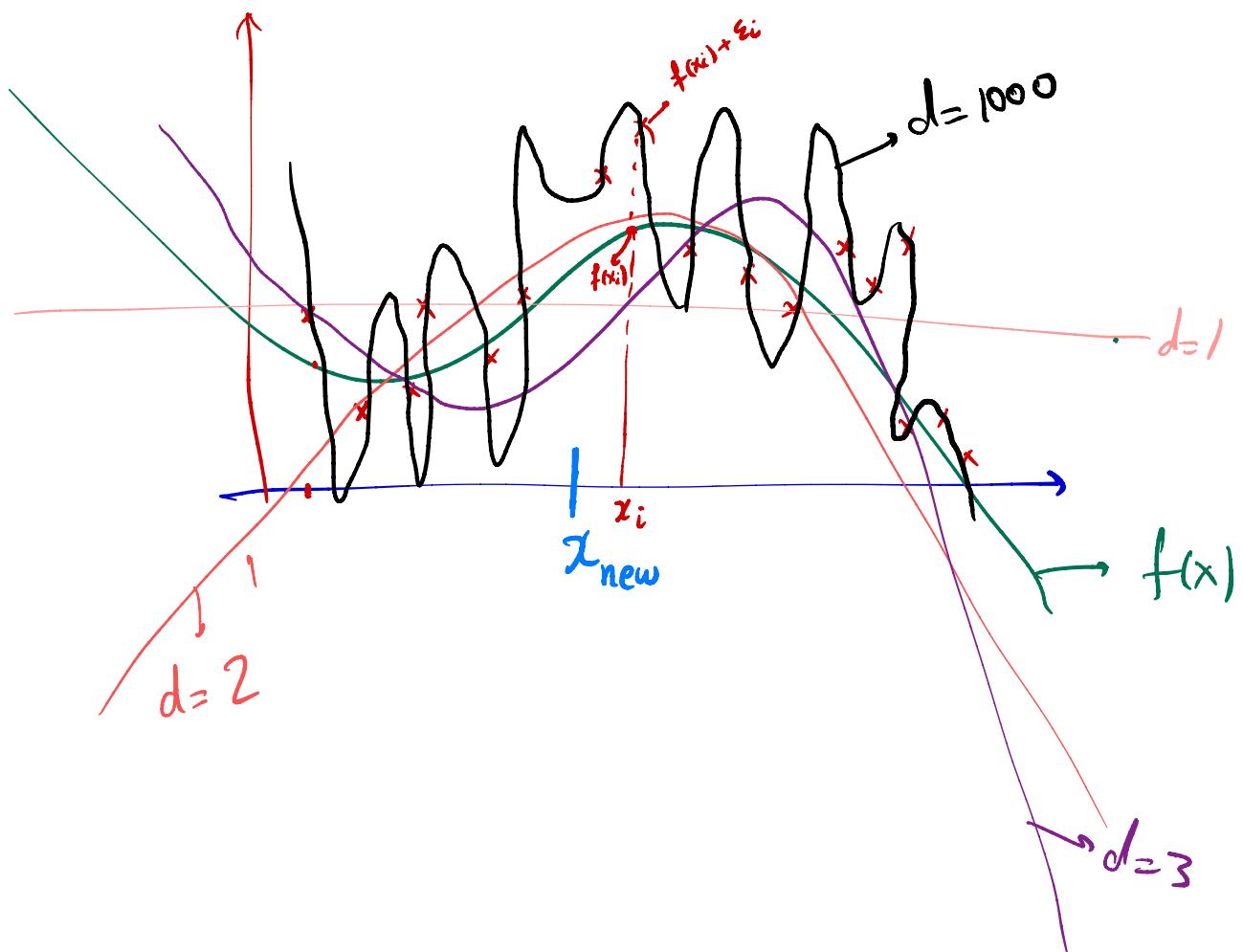
- Consider the following model for data:

how the data is generated

$$y = f(x) + \varepsilon$$

noise

e.g.  $f(x) = x^3 + 3x^2 - 4x + 5$



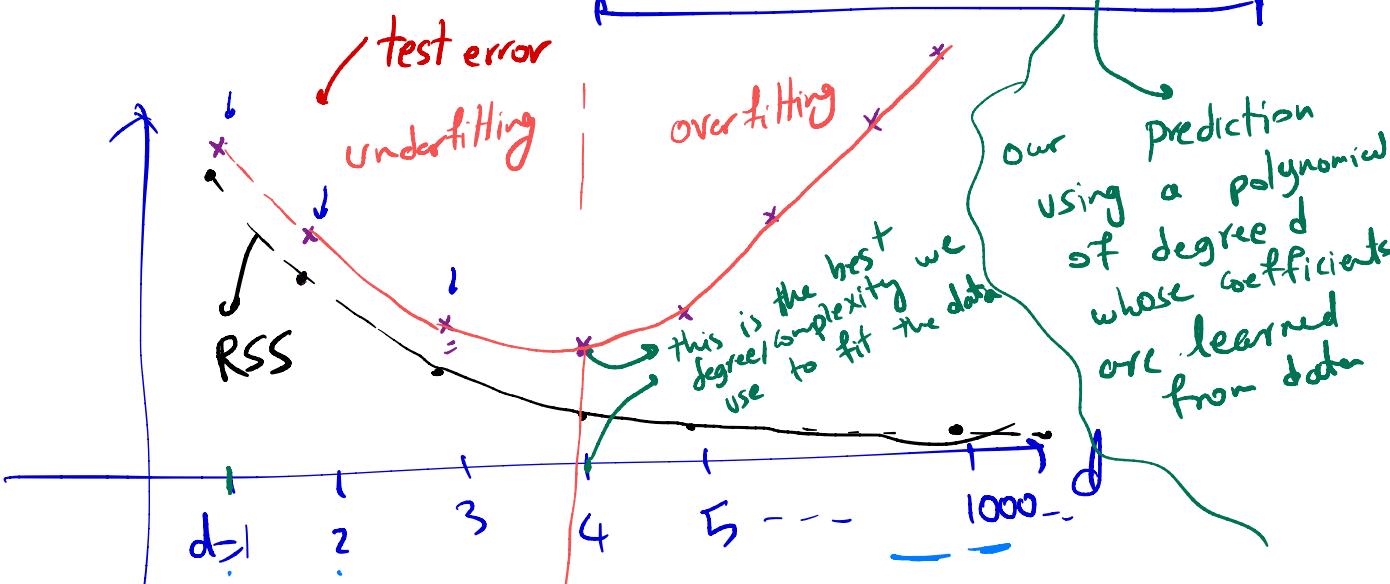
- On a new data point  $x_{\text{new}}$

$$\rightarrow y_{\text{new}} = \underline{f(x_{\text{new}})} + \varepsilon$$

w

the model learned from data

$$\rightarrow \text{Test error} = E \left[ (y_{\text{new}} - \underbrace{P_d(x_{\text{new}}, \vec{\beta})}_{y_{\text{new}} = f(x_{\text{new}}) + \varepsilon})^2 \right]$$



we have two regions depending on the complexity (degree) of the parametric model that we're fitting to the data:

- Underfitting:
  - the model is not sufficiently complex to explain the data (high bias)
  - the trained model does not change much if we add new training data points (low variance)
- Overfitting:
  - model is too complex for training the data (low bias)
  - the trained model will change significantly if we add new data points (high variance)