#### CS 453X: Class 6

Jacob Whitehill

# Gradient descent (continued)

#### Gradient descent algorithm

- Set w to random values; call this initial choice w<sup>(0)</sup>.
- Compute the gradient:  $\nabla_{\mathbf{w}} f(\mathbf{w}^{(0)})$
- Update **w** by moving opposite the gradient, multiplied by a step size  $\varepsilon$ .  $\mathbf{w}^{(1)} \leftarrow \mathbf{w}^{(0)} \epsilon \nabla_{\mathbf{w}} f(\mathbf{w}^{(0)})$
- Repeat...

$$\mathbf{w}^{(2)} \leftarrow \mathbf{w}^{(1)} - \epsilon \nabla_{\mathbf{w}} f(\mathbf{w}^{(1)})$$

$$\mathbf{w}^{(3)} \leftarrow \mathbf{w}^{(2)} - \epsilon \nabla_{\mathbf{w}} f(\mathbf{w}^{(2)})$$

...

$$\mathbf{w}^{(t)} \leftarrow \mathbf{w}^{(t-1)} - \epsilon \nabla_{\mathbf{w}} f(\mathbf{w}^{(t-1)})$$

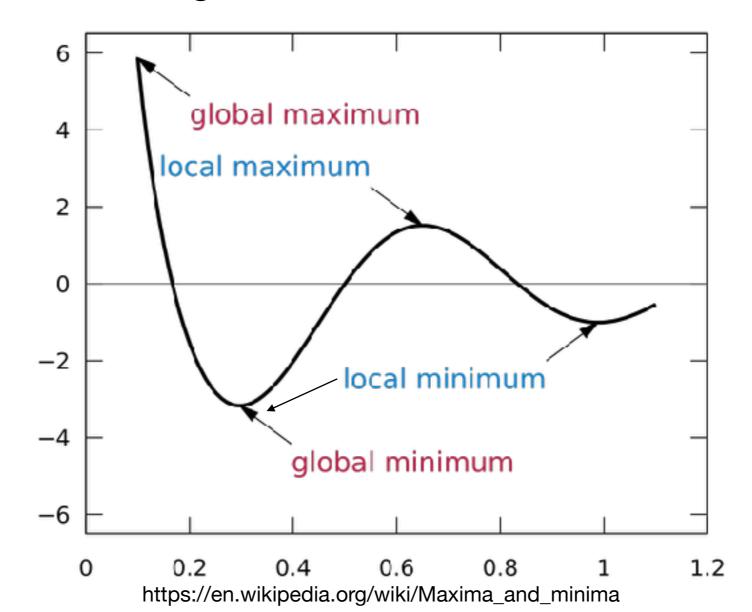
...until convergence:

$$|f(\mathbf{w}^{(t-1)}) - f(\mathbf{w}^{(t)})| < \delta \quad \text{ois a chosen convergence tolerance}$$

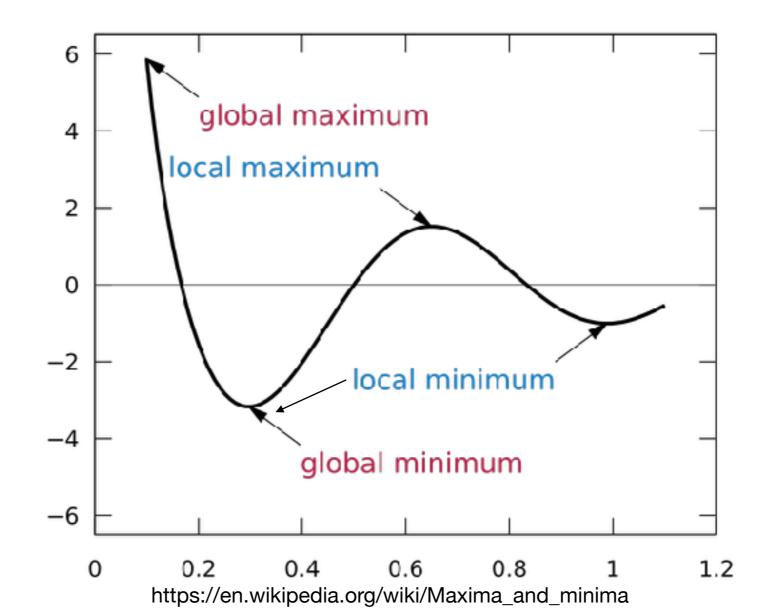
#### Gradient descent demos

- 1-d
- 2-d

- In general, gradient descent is useful for finding a local minimum of f:
  - Local minimum: gradient is 0; second derivative is positive.



- In general, gradient descent is useful for finding a local minimum of f:
  - Global minimum is the smallest value of the function f.



• For the special case of linear regression (and a few other ML models), gradient descent will (for appropriate  $\varepsilon$ ) converge to the *global* minimum of  $f_{MSE}$ .

- For the special case of linear regression (and a few other ML models), gradient descent will (for appropriate  $\varepsilon$ ) converge to the *global* minimum of  $f_{MSE}$ .
- What does "appropriate ε" mean (intuitively)?
  - Big enough to make progress (from random starting point) to local minimum.
  - Small enough not to "jump around" too much.
    - Show demo.

- For the special case of linear regression (and a few other ML models), gradient descent will (for appropriate  $\varepsilon$ ) converge to the *global* minimum of  $f_{MSE}$ .
- In practice:
  - Choose some  $\varepsilon$  so that the cost  $f_{MSE}$  declines smoothly.
  - If it's too slow, try increasing.
  - If it jumps around, try decreasing.

#### Linear regression

- Linear regression is efficient to optimize and very useful, but is limited in its expressiveness.
- Unsurprisingly, it can only model linear (technically affine) relationships:

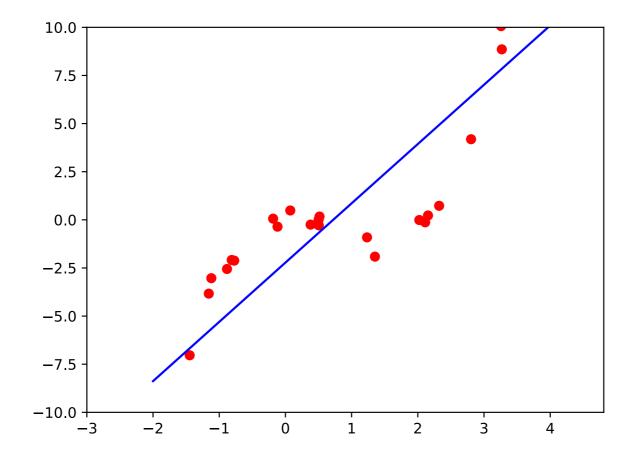
$$\hat{y} = \mathbf{x}^{\top} \mathbf{w} + b$$

• In 1-d:

$$\hat{y} = wx + b$$

#### Linear regression

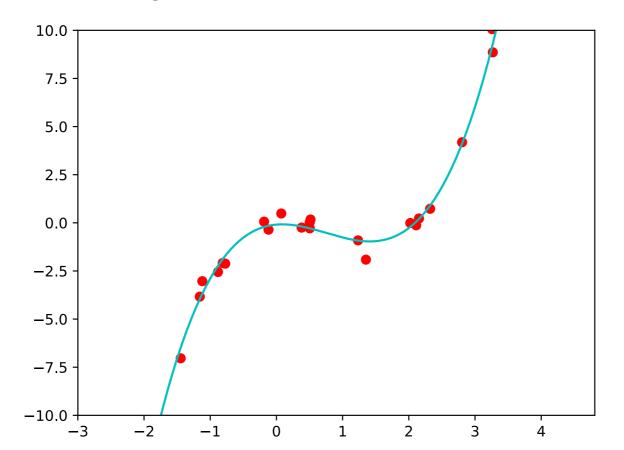
- But sometimes the target values y have a non-linear relationship with the input x.
- Linear regression may not do a good job then.
- Example:  $y = 0.2x 1.8x^2 + 0.8x^3 + \text{noise}$



Not a great fit.

- If the labels y are a polynomial function of the inputs x, why
  not enable the model to express polynomial relationships?
- In 1-d, we can build a cubic regression model as follows:

$$\hat{y} = w_1 x^1 + w_2 x^2 + w_3 x^3 + b$$
$$= w_0 x^0 + w_1 x^1 + w_2 x^2 + w_3 x^3$$



**Much better fit!** 

- How do we train the weights of the polynomial regression (for 1-d inputs)?
  - Pretend that each power of x is a separate feature.
  - Form  $\mathbf{x} = [x^0, x^1, x^2, x^3]^T$

- How do we train the weights of the polynomial regression (for 1-d inputs)?
  - Pretend that each power of x is a separate feature.
  - Form  $\mathbf{x} = [x^0, x^1, x^2, x^3]^T$
- Example:
  - Suppose the raw input x = -1.5.
  - Then  $x_0 = 1$ ,  $x_1 = -1.5$ ,  $x^2 = 2.25$ , and  $x^3 = -3.375$ . Hence,  $\mathbf{x} = [1, -1.5, 2.25, -3.375]^T$ .

Now, notice that:

$$\hat{y} = w_0 x^0 + w_1 x^1 + w_2 x^2 + w_3 x^3$$

Now, notice that:

$$\hat{y} = w_0 x^0 + w_1 x^1 + w_2 x^2 + w_3 x^3 
= \begin{bmatrix} x^0 & x^1 & x^2 & x^3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

Now, notice that:

$$\hat{y} = w_0 x^0 + w_1 x^1 + w_2 x^2 + w_3 x^3$$

$$= \begin{bmatrix} x^0 & x^1 & x^2 & x^3 \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

$$= \mathbf{x}^{\top} \mathbf{w}$$

When we "pre-compute" each power of x, we convert the polynomial regression back into a linear regression model.

 We can convert each input x into a feature vector x, and create a design matrix X, and compute the optimal w as before...

Suppose we have raw inputs -1.5, -1, and 3.25.

- Suppose we have raw inputs -1.5, -1, and 3.25.
- Then for each (scalar) x we build a vector x consisting of [x<sup>0</sup>, x<sup>1</sup>, x<sup>2</sup>, x<sup>3</sup>]<sup>T</sup>:

|  |             | <i>i</i> =1 | i=2 | <i>i</i> =3 |
|--|-------------|-------------|-----|-------------|
|  | <i>d</i> =0 | 1           | 1   | 1           |
|  | <i>d</i> =1 | -1.5        | -1  | 3.25        |
|  | d=2         | 2.25        | 1   | 10.5625     |
|  | d=3         | -3.375      | -1  | 4.32812     |

 The matrix of all our examples (as column vectors) constitutes the design matrix X, as usual.

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 \\ -1.5 & -1 & 3.25 \\ 2.25 & 1 & 10.5625 \\ -3.375 & -1 & 4.32812 \end{bmatrix}$$

- The matrix of all our examples (as column vectors) constitutes the design matrix X, as usual.
- We can now find the optimal polynomial regression coefficients by computing:

$$\mathbf{w} = \left(\mathbf{X}\mathbf{X}^{\top}\right)^{-1}\mathbf{X}\mathbf{y}$$

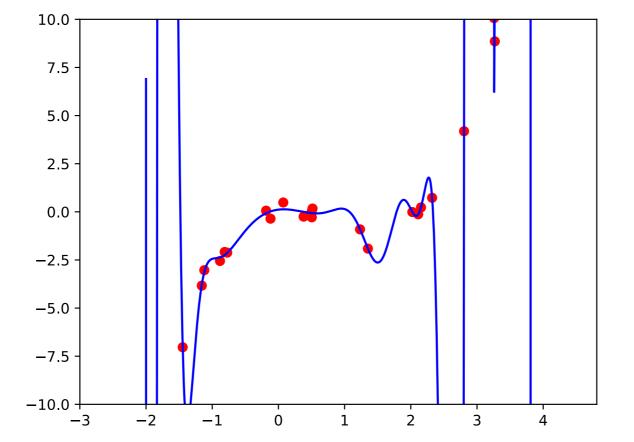
...just like with linear regression.

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 1 \\ -1.5 & -1 & 3.25 \\ 2.25 & 1 & 10.5625 \\ -3.375 & -1 & 4.32812 \end{bmatrix}$$

# Overfitting and regularization

- If polynomial regression with degree 3 worked well, why not increase the degree even higher?
- Let's try with degree 25...

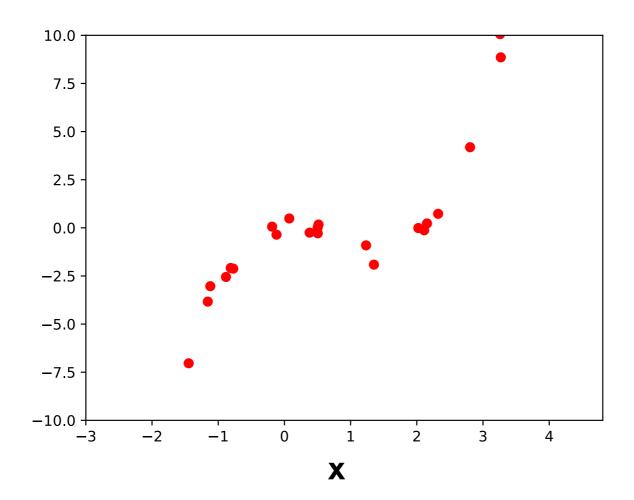
- If polynomial regression with degree 3 worked well, why not increase the degree even higher?
- Let's try with degree 25...



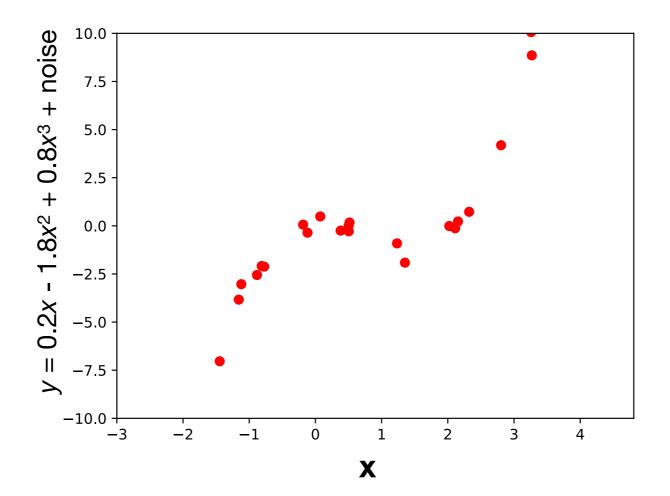
We nailed almost every point exactly!... but maybe this is overkill?

- Why is this bad? Recall that **overfitting** means that training error is low, but testing error is high.
- Testing error represents how well we expect our machine to perform on data we have not seen before.

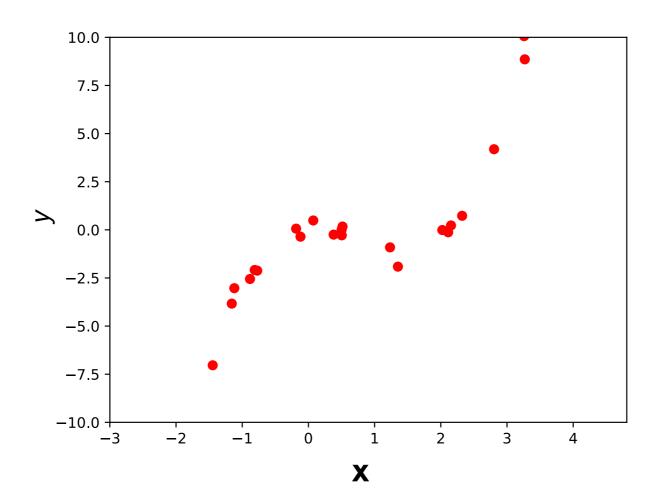
• When we collect a dataset, we are sampling from probability distribution  $p(\mathbf{x})$ 



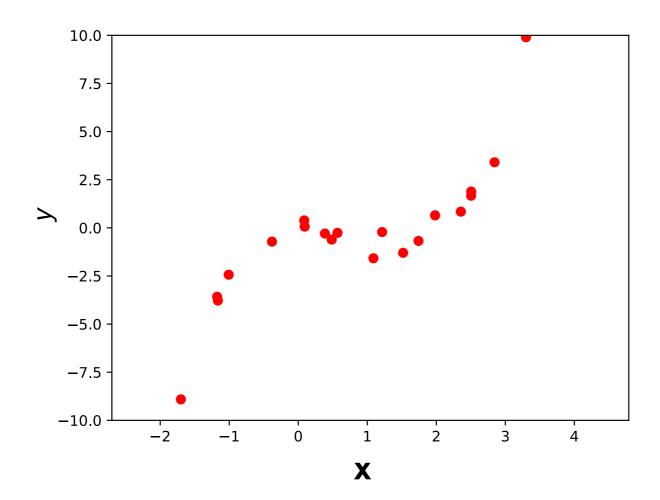
• When we collect a dataset, we are sampling from probability distribution  $p(\mathbf{x})$  and conditional probability distribution  $p(\mathbf{y} \mid \mathbf{x})$ .



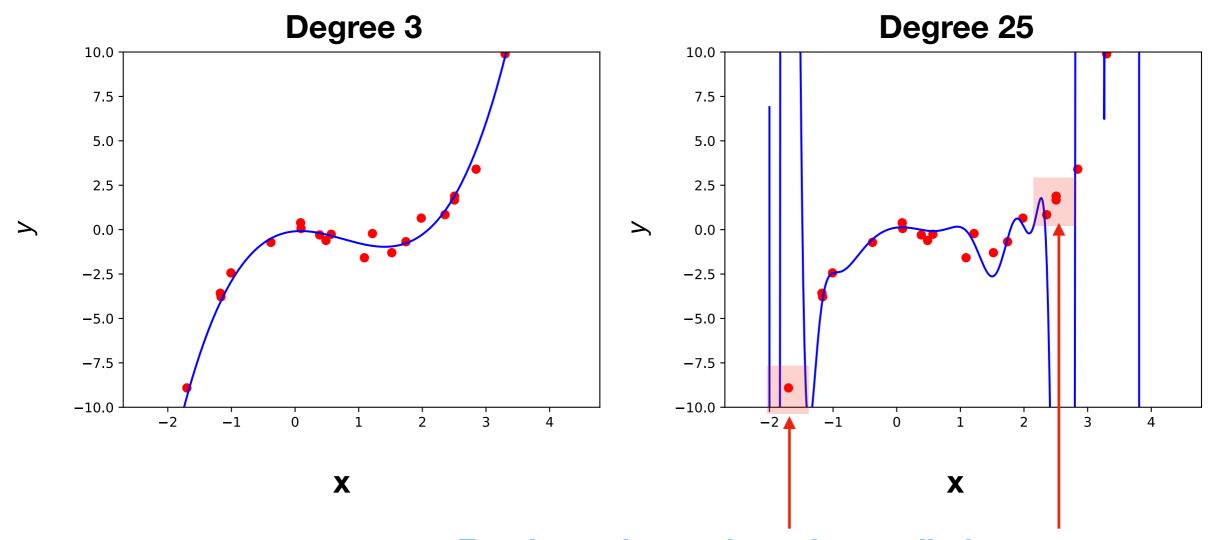
• When we sample multiple times, we will get different results. Here is a possible *training* sample:



• When we sample multiple times, we will get different results. Here is a possible *testing* sample:



 Here are the machine's predictions using polynomial regression, with either degree 3 or degree 5:



For these data points, the predictions are very inaccurate, which makes  $f_{MSE}$  large.

#### Preventing overfitting

- How to prevent this? Two strategies:
  - Keep the degree d of the polynomial modest.

#### Preventing overfitting

- How to prevent this? Two strategies:
  - Keep the degree d of the polynomial modest.
  - Keep the weight associated with each term modest.

$$\hat{y} = w_0 x^0 + w_1 x^1 + w_2 x^2 + \ldots + w_d x^d$$

 Let's generate some polynomials by randomly selecting each w<sub>i</sub> in:

$$\hat{y} = w_0 x^0 + w_1 x^1 + w_2 x^2 + \ldots + w_d x^d$$

Compute the average squared coefficient as:

$$\mu = \frac{1}{d} \sum_{i=0}^{d} w_i^2$$

We will generate random polynomials for different degrees
 d and different coefficient magnitudes μ.

$$\mu = \frac{1}{d} \sum_{i=0}^{d} w_i^2$$

Examples:

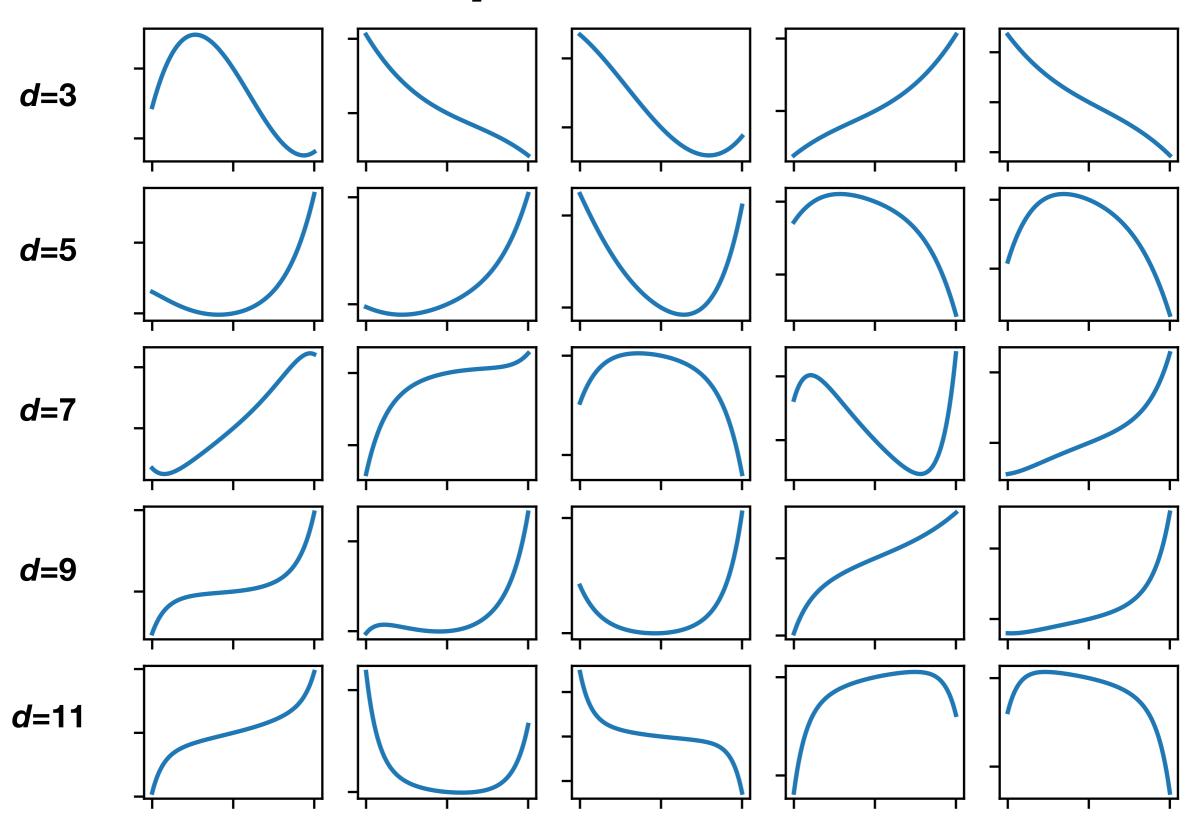
$$d=2: \quad \hat{y}=4+2x-2x^2 \qquad \mu = 2$$
 $d=4: \quad \hat{y}=x^2+0.5x^3-2x^4 \quad \mu = 2$ 

$$\mu = \frac{1}{d} \sum_{i=0}^{d} w_i^2$$

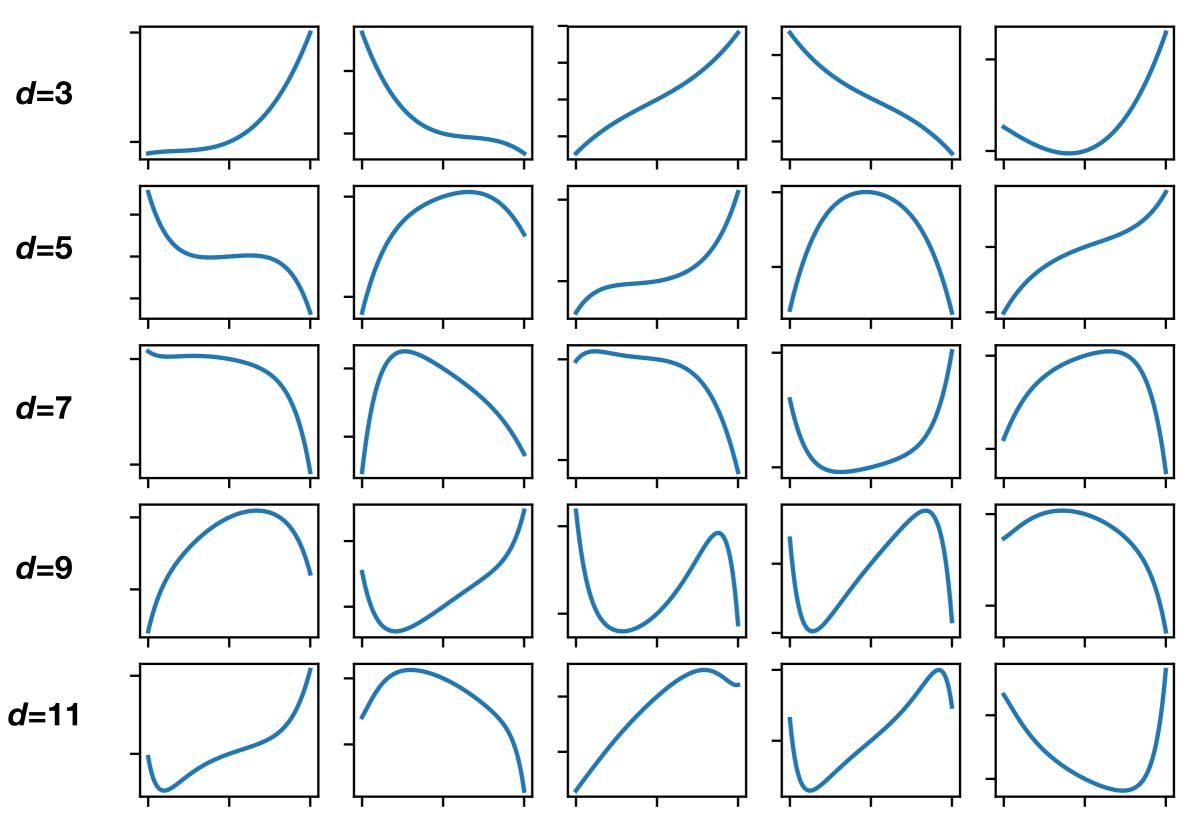
#### Examples:

$$d = 2: \quad \hat{y} = 4 + 2x - 2x^2 \qquad \qquad \mu = 24/3 = 8$$
 
$$d = 4: \quad \hat{y} = x^2 + 0.5x^3 - 2x^4 \qquad \mu = 5.25/5 = 1.05$$

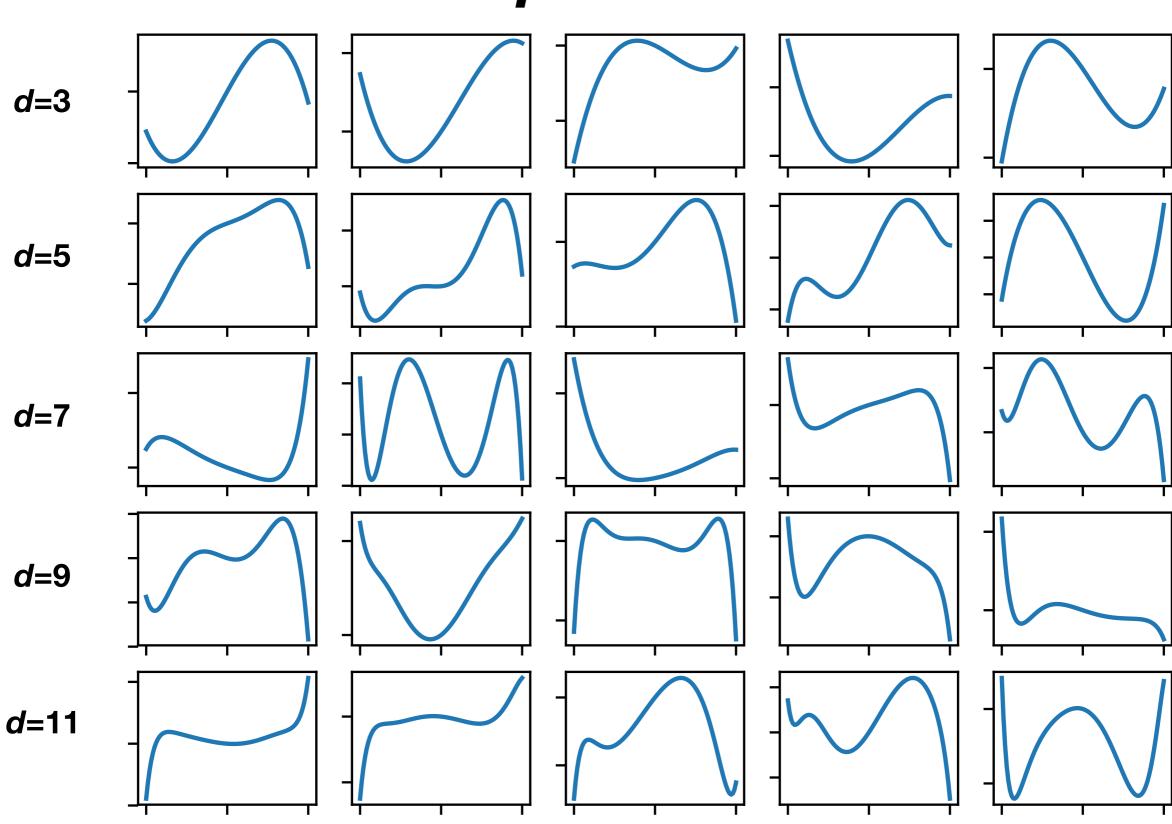
 $\mu$ =7.25e-07



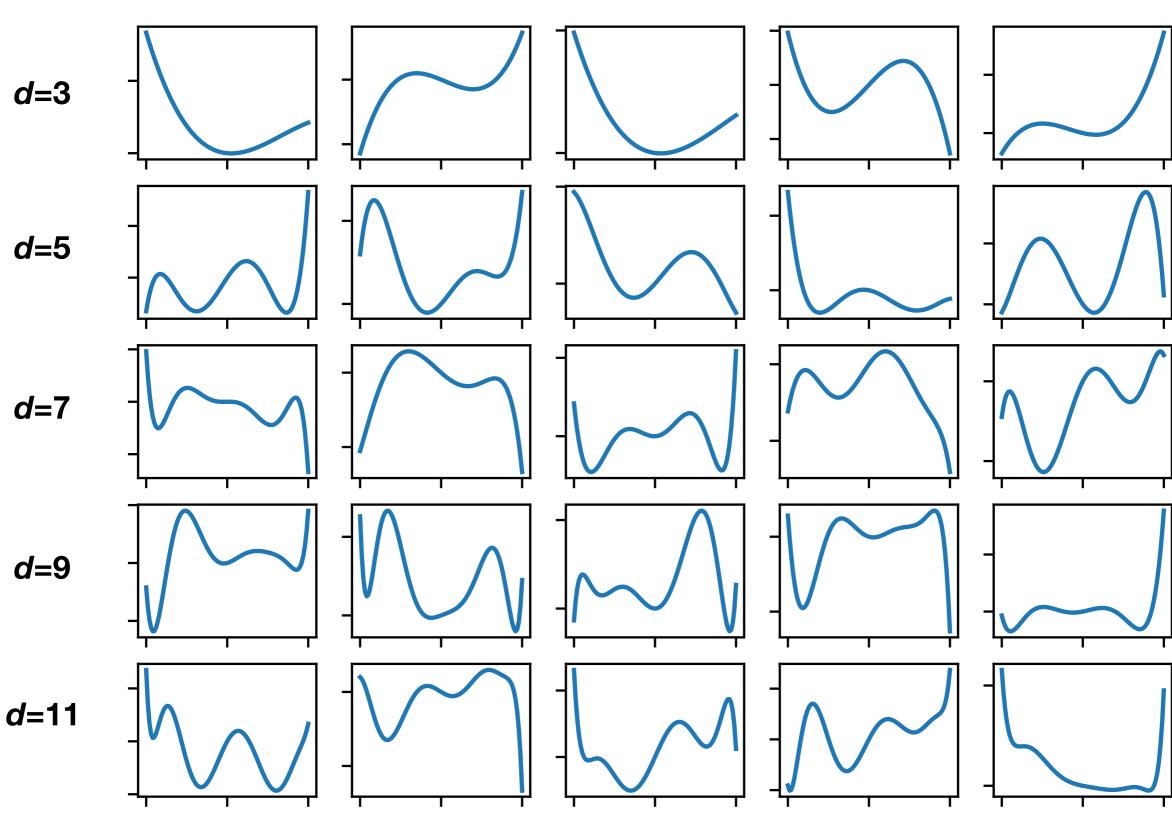
 $\mu$ =0.00063



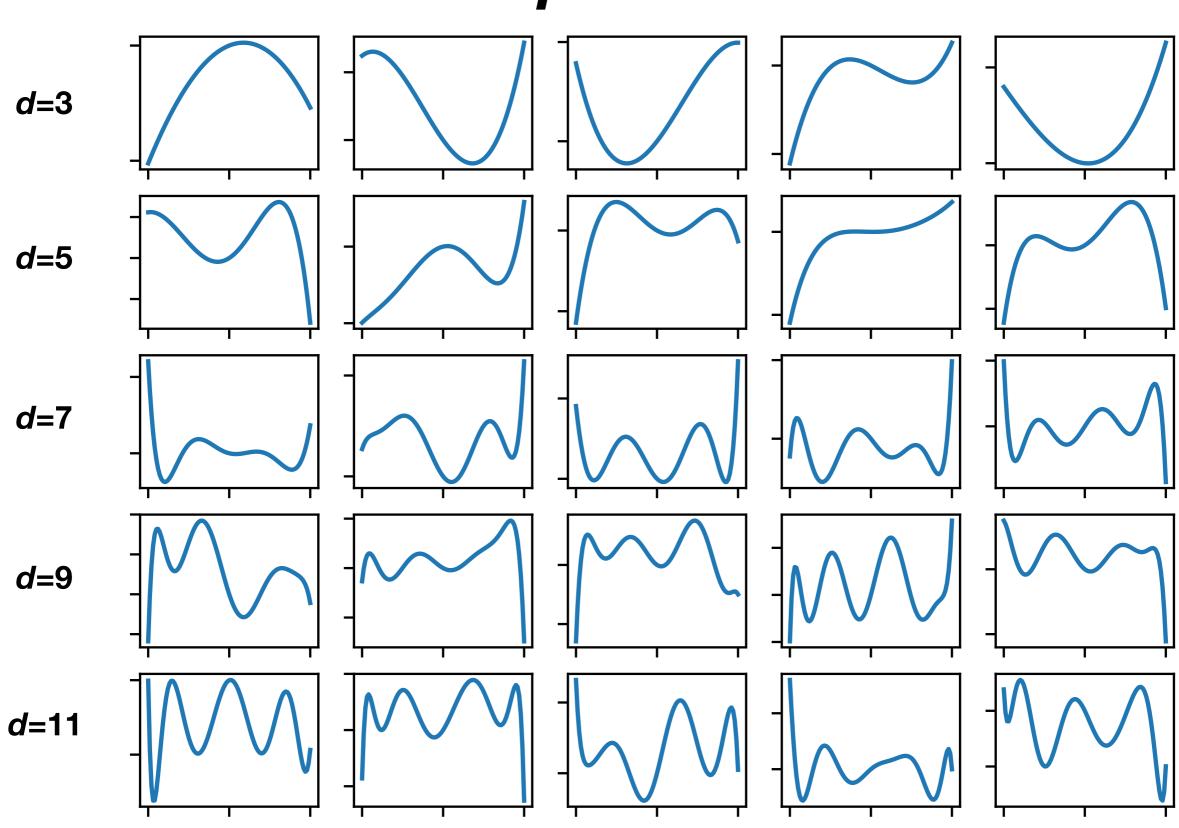
 $\mu$ =0.058



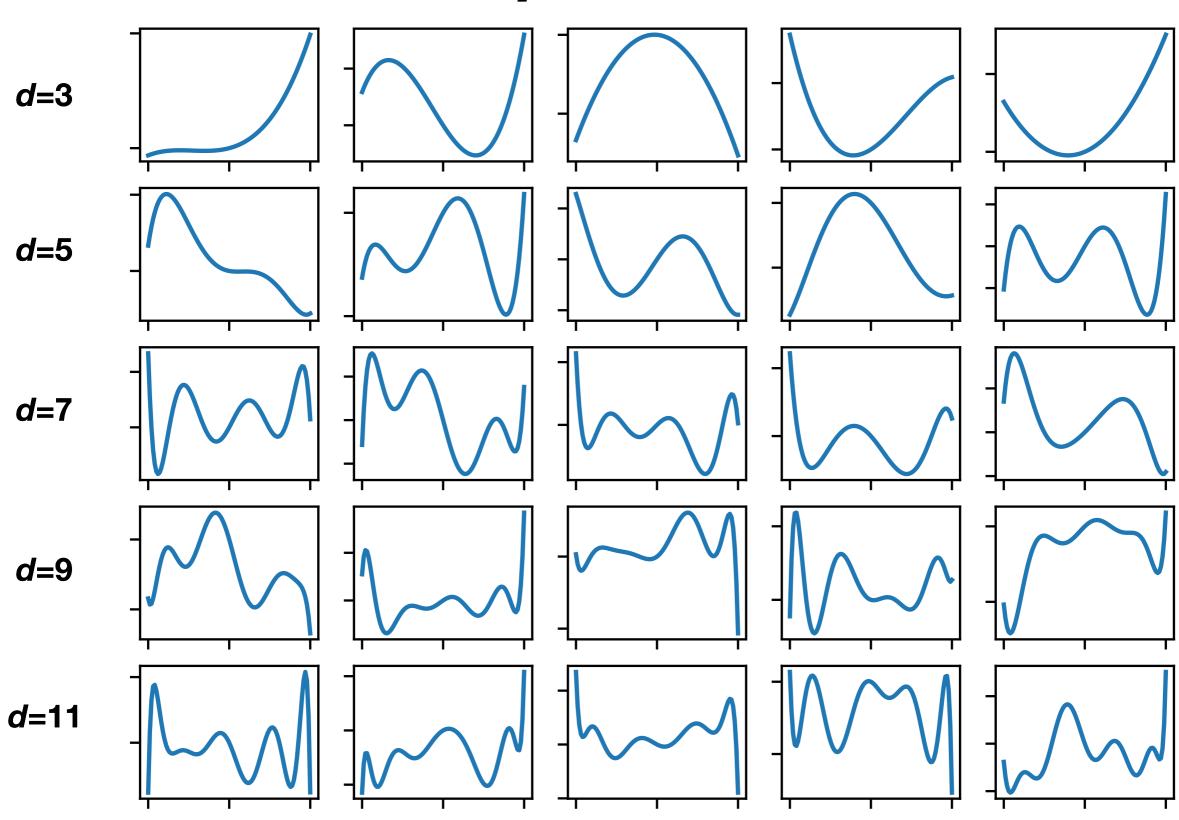
 $\mu$ =3.31



 $\mu = 90.9$ 



 $\mu = 537.8$ 



#### Regularization

- The larger the coefficients (weights) **w** are allowed to be, the more the polynomial regressor can overfit.
- If we "encourage" the weights to be small, we can reduce overfitting.
- This is a form of regularization any practice designed to improve the machine's ability to generalize to new data.

#### Regularization

- One of the simplest and oldest regularization techniques is to penalize large weights in the cost function.
- The "unregularized"  $f_{MSE}$  is:

$$f_{\text{MSE}}(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^{n} (y^{(i)} - \hat{y}^{(i)})^2$$

#### Regularization

- One of the simplest and oldest regularization techniques is to penalize large weights in the cost function.
- The "unregularized"  $f_{MSE}$  is:

$$f_{\text{MSE}}(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^{n} (y^{(i)} - \hat{y}^{(i)})^2$$

• The *L*<sub>2</sub>-regularized *f*<sub>MSE</sub> becomes:

$$f_{\text{MSE}}(\mathbf{w}) = \frac{1}{2n} \sum_{i=1}^{n} (y^{(i)} - \hat{y}^{(i)})^2 + \frac{\alpha}{2} \mathbf{w}^{\top} \mathbf{w}$$