



# 统计机器学习 (小班研讨)

主讲教师：刘峤

# 回归算法

An introduction to regression

# Three interpretations of regression

- Linear regression

$$\hat{y} = \mathbf{w} \cdot \mathbf{x}$$

- Probabilistic (Maximum Likelihood Estimation, MLE)

$$y \sim \mathcal{N}(\mathbf{w} \cdot \mathbf{x}, \sigma^2)$$

$$\underset{\mathbf{w}}{\operatorname{argmax}} p(y|\mathbf{w}, \mathbf{x})$$

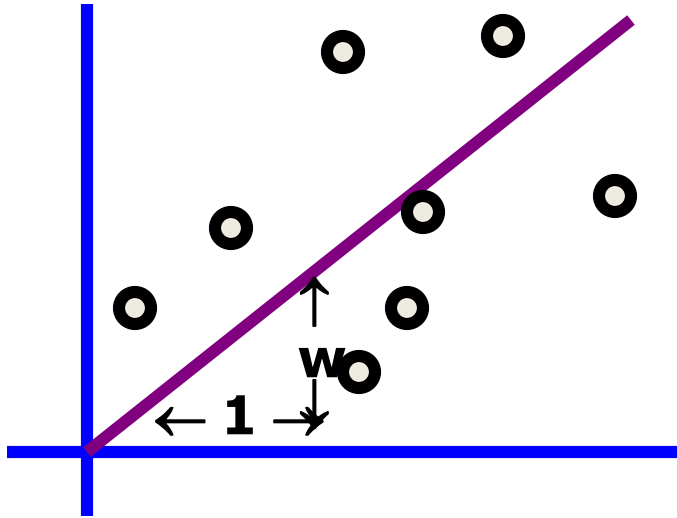
- Bayesian (Maximum Posterior Probability Estimation, MAP)

$$\underset{\mathbf{w}}{\operatorname{argmax}} p(y|\mathbf{w}, \mathbf{x})p(\mathbf{w}, \mathbf{x})$$

# Single-Parameter Linear Regression

# Linear Regression

**DATASET**



inputs	outputs
$x_1 = 1$	$y_1 = 1$
$x_2 = 3$	$y_2 = 2.2$
$x_3 = 2$	$y_3 = 2$
$x_4 = 1.5$	$y_4 = 1.9$
$x_5 = 4$	$y_5 = 3.1$

- Linear regression assumes that the expected value of the output given an input,  $E[y/x]$ , is linear.
- Simplest case:  $\text{Out}(x) = wx$  for some unknown  $w$ .
- Given the data, we can estimate  $w$ .

# 1-parameter linear regression

Assume that the data is formed by

$$y_i = wx_i + \text{noise}_i$$

where...

- the noise signals are independent
- the noise has a normal distribution with mean 0 and unknown variance  $\sigma^2$

Means?

- $p(y | w, x)$  has a normal distribution with
  - mean  $wx$
  - variance  $\sigma^2$

# Bayesian Linear Regression

$$p(y|w,x) = \text{Normal}(\text{mean } wx, \text{var } \sigma^2)$$

- We have a set of datapoints  $(x_1, y_1)$   $(x_2, y_2)$  ...  $(x_n, y_n)$  which are **EVIDENCE** about  $w$ .

- We want to infer  $w$  from the data.

$$p(w|x_1, x_2, x_3, \dots, x_n, y_1, y_2, \dots, y_n)$$

- You can use **BAYES** rule to work out a posterior distribution for  $w$  given the data.
- Or you could do Maximum Likelihood Estimation

# Maximum likelihood estimation of $w$

MLE asks: For which value of  $w$  is this data most likely to have happened?

$\Leftrightarrow$  For what  $w$  is

$p(y_1, y_2 \dots y_n \mid x_1, x_2, x_3, \dots x_n, w)$  maximized?

$\Leftrightarrow$  For what  $w$  is

$$\prod_{i=1}^n p(y_i \mid w, x_i) \text{ maximized?}$$



For what  $w$  is

$$\prod_{i=1}^n p(y_i | w, x_i) \text{ maximized?}$$

For what  $w$  is

$$\prod_{i=1}^n \exp\left(-\frac{1}{2} \left(\frac{y_i - wx_i}{\sigma}\right)^2\right) \text{ maximized?}$$

For what  $w$  is

$$\sum_{i=1}^n -\frac{1}{2} \left(\frac{y_i - wx_i}{\sigma}\right)^2 \text{ maximized?}$$

For what  $w$  is

$$\sum_{i=1}^n (y_i - wx_i)^2 \text{ minimized?}$$

# First result

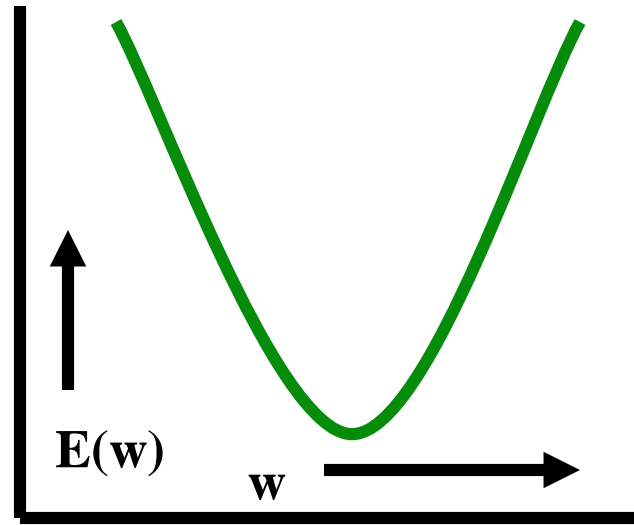
---

- MLE with **Gaussian noise** is the same as minimizing the  $L_2$  error

$$\arg \min \sum_{i=1}^n (y_i - w x_i)^2$$

# Linear Regression

The maximum likelihood  $w$  is the one that minimizes sum-of-squares of residuals



$$\begin{aligned} E &= \sum_i (y_i - wx_i)^2 \\ &= \sum_i y_i^2 - (2 \sum x_i y_i)w + \left( \sum x_i^2 \right) w^2 \end{aligned}$$

We want to minimize a quadratic function of  $w$ .

# Linear Regression

**Easy to show the sum of squares is minimized when**

$$w = \frac{\sum x_i y_i}{\sum x_i^2}$$

**The maximum likelihood model is**

$$\text{Out}(x) = wx$$

**We can use it for prediction**

# But what about MAP?

- **MLE**

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n p(y_i | \mathbf{w}, x_i)$$

- **MAP**

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n p(y_i | \mathbf{w}, x_i) p(\mathbf{w})$$

# But what about MAP?

- MAP

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n p(y_i | \mathbf{w}, x_i) p(\mathbf{w})$$

- We assumed

$$y \sim \mathcal{N}(\mathbf{w} \cdot \mathbf{x}, \sigma^2)$$

- Now add a prior that assumption that

$$\mathbf{w} \sim \mathcal{N}(0, \gamma^2)$$

For what  $w$  is

$$\prod_{i=1}^n p(y_i | \mathbf{w}, x_i) p(\mathbf{w}) \quad \text{maximized?}$$

For what  $w$  is

$$\prod_{i=1}^n \exp \left( -\frac{1}{2} \left( \frac{y_i - \mathbf{w}x_i}{\sigma} \right)^2 \right) \exp \left( -\frac{1}{2} \left( \frac{\mathbf{w}}{\gamma} \right)^2 \right) \quad \text{maximized?}$$

For what  $w$  is

$$\sum_{i=1}^n -\frac{1}{2} \left( \frac{y_i - \mathbf{w}x_i}{\sigma} \right)^2 - \frac{1}{2} \left( \frac{\mathbf{w}}{\gamma} \right)^2 \quad \text{maximized?}$$

For what  $w$  is

$$\sum_{i=1}^n (y_i - \mathbf{w}x_i)^2 + \left( \frac{\sigma \mathbf{w}}{\gamma} \right)^2 \quad \text{minimized?}$$

# Second result

- MAP with **Gaussian prior** on  $w$  is the same as minimizing the  $L_2$  error plus an  $L_2$  penalty on  $w$

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w}x_i)^2 + \lambda \mathbf{w}^2$$

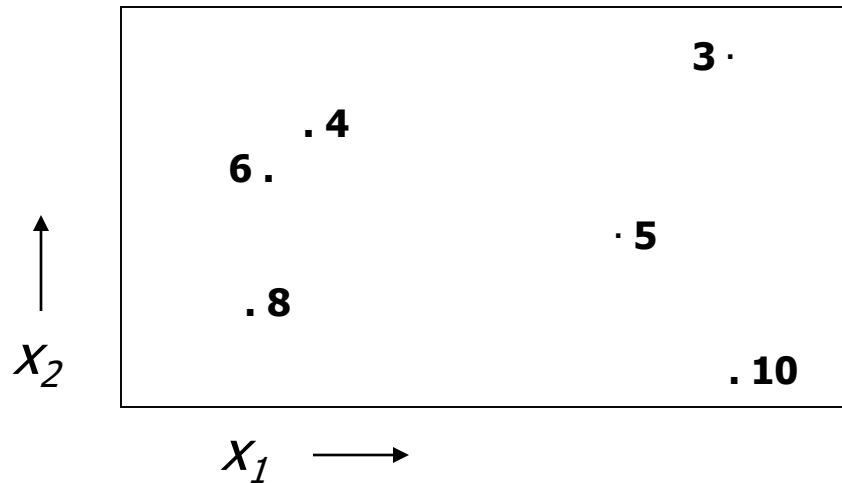
- This is called: **Regularization**
  - **Ridge regression**
  - **Shrinkage**



# Multivariate Linear Regression

# Multivariate Regression

What if the inputs are vectors?



**2-d input  
example**

Dataset has form

$$\begin{array}{cc} \mathbf{x}_1 & y_1 \\ \mathbf{x}_2 & y_2 \\ \mathbf{x}_3 & y_3 \\ \vdots & \vdots \\ \mathbf{x}_R & y_R \end{array}$$

# Multivariate Regression

Write matrix  $X$  and  $Y$  thus:

$$\mathbf{X} = \begin{bmatrix} \dots \mathbf{X}_1 \dots \\ \dots \mathbf{X}_2 \dots \\ \vdots \\ \dots \mathbf{X}_R \dots \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & x_{22} & \dots & x_{2m} \\ & & \vdots & \\ x_{R1} & x_{R2} & \dots & x_{Rm} \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_R \end{bmatrix}$$

(there are  $R$  datapoints. Each input has  $m$  components)

The linear regression model assumes a vector  $\mathbf{w}$  such that

$$\text{Out}(\mathbf{x}) = \mathbf{x} \mathbf{w} = w_1 x[1] + w_2 x[2] + \dots w_m x[D]$$

The max. likelihood  $\mathbf{w}$  is  $\mathbf{w} = (X^T X)^{-1} (X^T Y)$

# Multivariate Regression (con't)

The max. likelihood  $\mathbf{w}$  is  $\mathbf{w} = (X^T X)^{-1} (X^T Y)$

$X^T X$  is an  $m \times m$  matrix:  $i, j$ 'th elt is  $\sum_{k=1}^R x_{ki} x_{kj}$

$X^T Y$  is an  $m$ -element vector:  $i$ 'th elt  $\sum_{k=1}^R x_{ki} y_k$

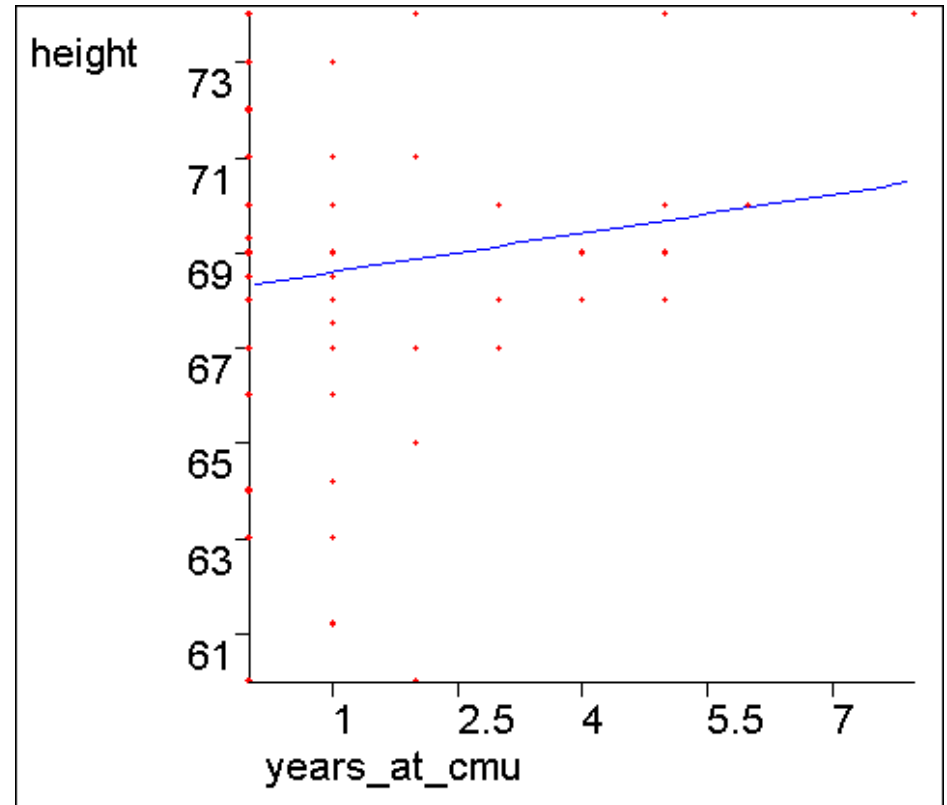
# Constant Term in Linear Regression

# What about a constant term?

We may expect linear data that does not go through the origin.

Statisticians and Neural Net Folks all agree on a simple obvious hack.

**Can you guess??**



# The constant term

- The trick is to create a fake input “ $X_0$ ” that always takes the value 1

$X_1$	$X_2$	$Y$
2	4	16
3	4	17
5	5	20

Before:

$$Y = w_1 X_1 + w_2 X_2$$

...has to be a poor model

$X_0$	$X_1$	$X_2$	$Y$
1	2	4	16
1	3	4	17
1	5	5	20

After:

$$Y = w_0 X_0 + w_1 X_1 + w_2 X_2$$
$$= w_0 + w_1 X_1 + w_2 X_2$$

...has a fine constant term

In this example,  
You should be  
able to see the  
MLE  $w_0$ ,  $w_1$  and  
 $w_2$  by inspection

Heteroscedasticity...

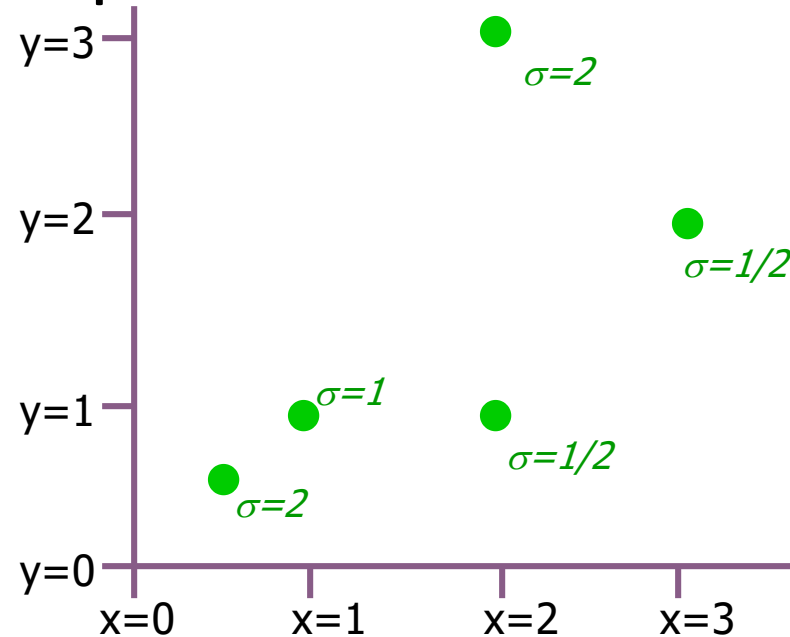
# Linear Regression with varying noise



# Regression with varying noise

- Suppose you know the variance of the noise that was added to each data point.

$x_i$	$y_i$	$\sigma_i^2$
$1/2$	$1/2$	4
1	1	1
2	1	$1/4$
2	3	4
3	2	$1/4$



Assume  $y_i \sim N(wx_i, \sigma_i^2)$

What's the MLE estimate of  $w$ ?

# MLE estimation with varying noise

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R \mid x_1, x_2, \dots, x_R, \sigma_1^2, \sigma_2^2, \dots, \sigma_R^2, w) =$$

$w$

$$\operatorname{argmin}_w \sum_{i=1}^R \frac{(y_i - wx_i)^2}{\sigma_i^2} =$$

Assuming independence among noise and then plugging in equation for Gaussian and simplifying.

$$\left( w \text{ such that } \sum_{i=1}^R \frac{x_i (y_i - wx_i)}{\sigma_i^2} = 0 \right) =$$

Setting  $dLL/dw$  equal to zero

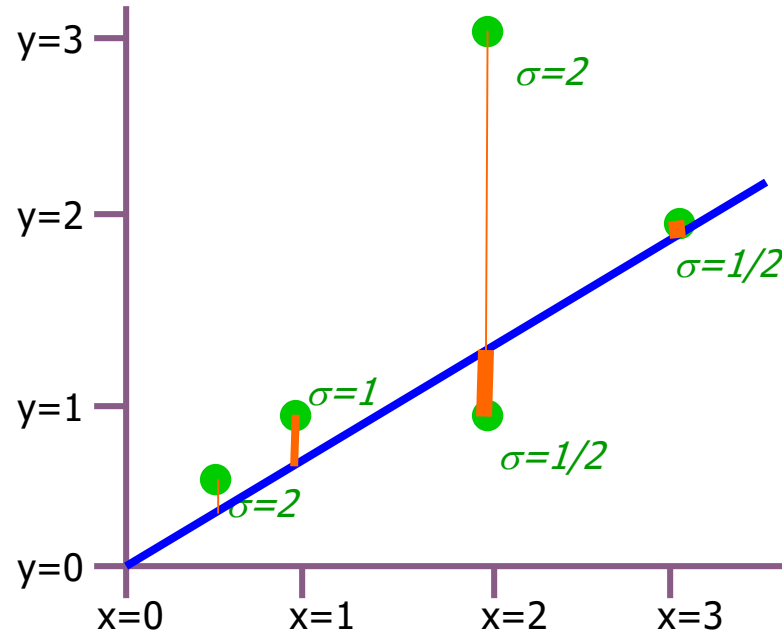
$$\frac{\left( \sum_{i=1}^R \frac{x_i y_i}{\sigma_i^2} \right)}{\left( \sum_{i=1}^R \frac{x_i^2}{\sigma_i^2} \right)}$$

Trivial algebra

# This is Weighted Regression

- We are asking to minimize the weighted sum of squares

$$\operatorname{argmin}_w \sum_{i=1}^R \frac{(y_i - wx_i)^2}{\sigma_i^2}$$



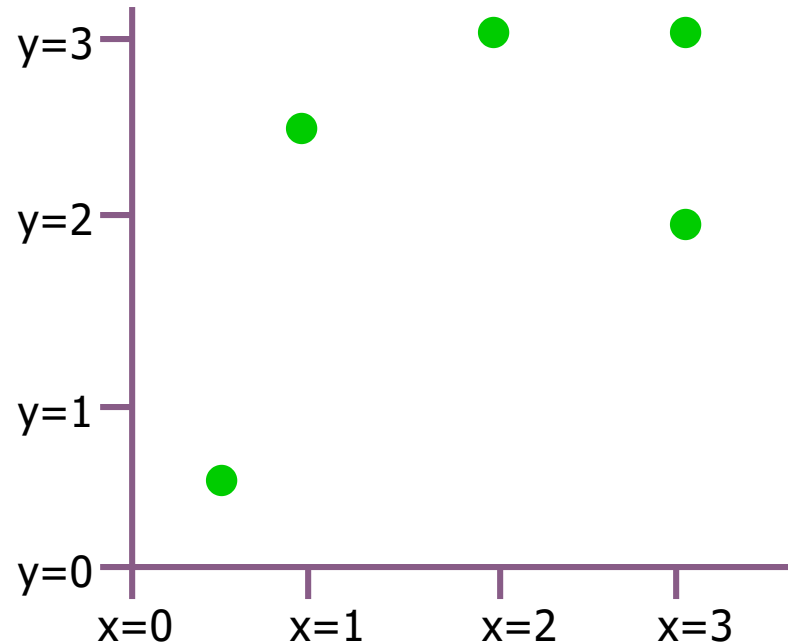
where weight for i'th datapoint is  $\frac{1}{\sigma_i^2}$

# Non-linear Regression

# Non-linear Regression

- Suppose you know that  $y$  is related to a function of  $x$  in such a way that the predicted values have a non-linear dependence on  $w$ , e.g:

$x_i$	$y_i$
$1/2$	$1/2$
1	2.5
2	3
3	2
3	3



Assume  $y_i \sim N(\sqrt{w + x_i}, \sigma^2)$

What's the MLE estimate of  $w$ ?

# Non-linear MLE estimation

$$\operatorname{argmax}_w \log p(y_1, y_2, \dots, y_R \mid x_1, x_2, \dots, x_R, \sigma, w) =$$

$$\operatorname{argmin}_w \sum_{i=1}^R \left( y_i - \sqrt{w + x_i} \right)^2 =$$

Assuming i.i.d. and then plugging in equation for Gaussian and simplifying.

$$\left( w \text{ such that } \sum_{i=1}^R \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Setting dLL/dw equal to zero

So guess what we do?

# Non-linear MLE estimation

$$\left( w \text{ such that } \sum_{i=1}^R \frac{y_i - \sqrt{w + x_i}}{\sqrt{w + x_i}} = 0 \right) =$$

Common (but not only) approach:

Numerical Solutions:

- Line Search
- **Simulated Annealing**
- **Gradient Descent**
- Conjugate Gradient
- Levenberg Marquart
- **Newton's Method**

Also, special purpose statistical-optimization-specific tricks such as E.M. (See Gaussian Mixtures lecture for introduction)

# Polynomial Regression



# Polynomial Regression

So far we've mainly been dealing with linear regression

$X_1$	$X_2$	$Y$
3	2	7
1	1	3
$\vdots$	$\vdots$	$\vdots$

$x =$

3	2
1	1
$\vdots$	$\vdots$

$y =$

7
3
$\vdots$

$y_1 = 7..$

$z =$

1	3	2
1	1	1
$\vdots$	$\vdots$	$\vdots$

$y =$

7
3
$\vdots$

$z_1 = (1, 3, 2)..$        $y_1 = 7..$

$z_k = (1, x_{k1}, x_{k2})$

$$\beta = (Z^T Z)^{-1} (Z^T y)$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

# Quadratic Regression

It's trivial to do linear fits of fixed nonlinear basis functions

$X_1$	$X_2$	$Y$
3	2	7
1	1	3
$\vdots$	$\vdots$	$\vdots$

$x =$		$y =$
3	2	7
1	1	3
$\vdots$	$\vdots$	$\vdots$

$y_1 = 7..$

$z =$						$y =$
1	3	2	9	6	4	7
1	1	1	1	1	1	3
$\vdots$						$\vdots$

$z = (1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$

$$\beta = (Z^T Z)^{-1} (Z^T y)$$

$$y^{est} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1^2 + \beta_4 x_1 x_2 + \beta_5 x_2^2$$

# Quadratic Regression

- $\mathbf{z}=(1, x_1, x_2, x_1^2, x_1x_2, x_2^2)$ 
  - Each component of a  $\mathbf{z}$  vector is called a term.
  - Each column of the  $\mathbf{Z}$  matrix is called a term column
- How many terms in a quadratic regression with  $m$  inputs?
  - 1 constant term
  - $m$  linear terms
  - $(m+1)\text{-choose-}2 = m(m+1)/2$  quadratic terms
  - $(m+2)\text{-choose-}2$  terms in total =  $O(m^2)$ .
- Note that solving  $\beta=(\mathbf{Z}^T\mathbf{Z})^{-1}(\mathbf{Z}^T\mathbf{y})$  is thus  $O(m^6)$

# Q<sup>th</sup>-degree polynomial Regression

$X_1$	$X_2$	$Y$
3	2	7
1	1	3
$\vdots$	$\vdots$	$\vdots$

$x=$	3	2	$y=$	7
	1	1		3
	$\vdots$	$\vdots$		$\vdots$

$z=$	1	3	2	9	6	...	$y=$	7
	1	1	1	1	1	...		3
	$\vdots$					...		$\vdots$

$z=(\text{all products of powers of inputs in which sum of powers is } q \text{ or less,})$

$$\beta = (Z^T Z)^{-1} (Z^T y)$$

$$y^{\text{est}} = \beta_0 + \beta_1 x_1 + \dots$$

m inputs, degree Q: how many terms?

= the number of unique terms of the form

$$x_1^{q_1} x_2^{q_2} \dots x_m^{q_m} \text{ where } \sum_{i=1}^m q_i \leq Q$$

= the number of unique terms of the form

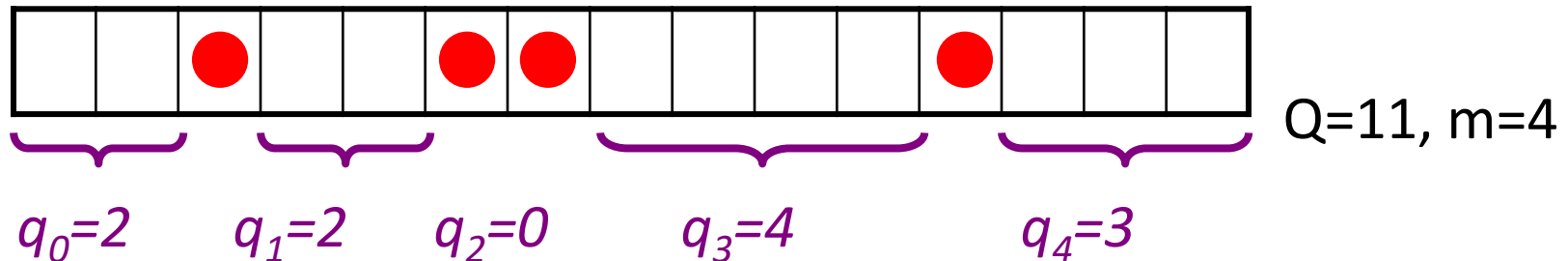
$$1^{q_0} x_1^{q_1} x_2^{q_2} \dots x_m^{q_m} \text{ where } \sum_{i=0}^m q_i = Q$$

= the number of lists of non-negative integers  $[q_0, q_1, q_2, \dots, q_m]$

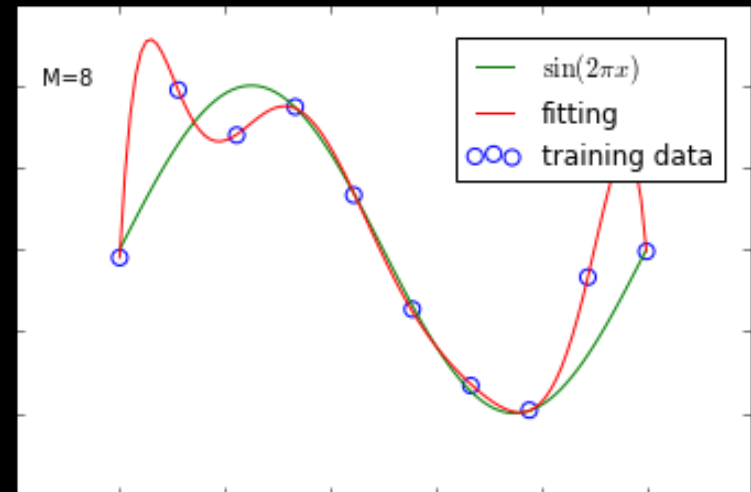
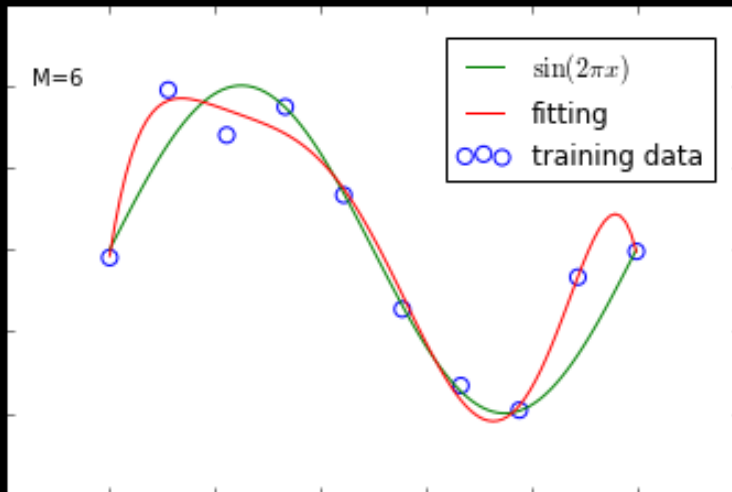
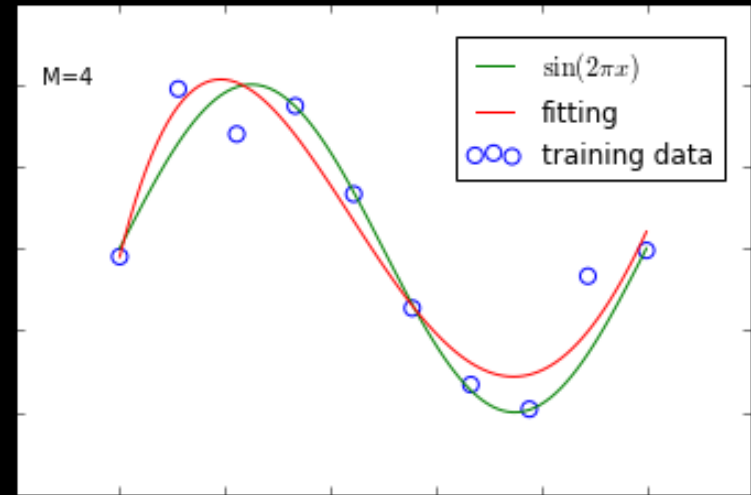
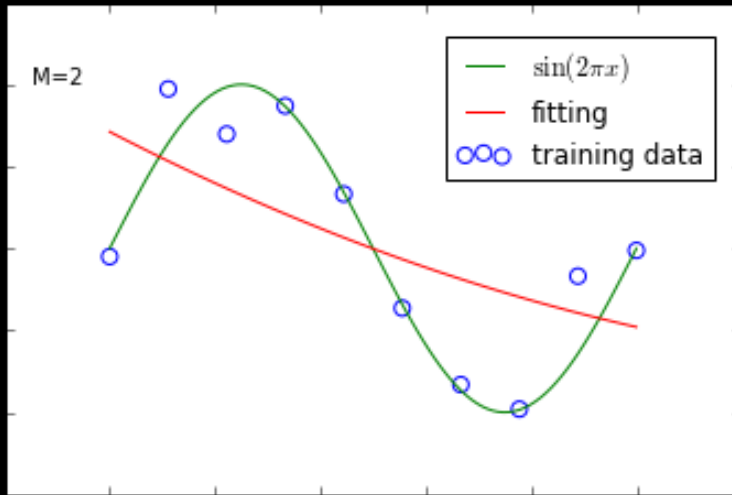
in which  $\sum q_i = Q$

= the number of ways of placing Q red disks on a row of squares of

length  $Q+m = (Q+m)\text{-choose-}Q$



# Polynomial Regression



# Multivariate vs Ridge Regression

- The linear regression model assumes a vector  $\mathbf{w}$  such that

$$\text{Out}(\mathbf{x}) = \mathbf{x} \mathbf{w} = w_1 x_1 + w_2 x_2 + \dots w_m x_m$$

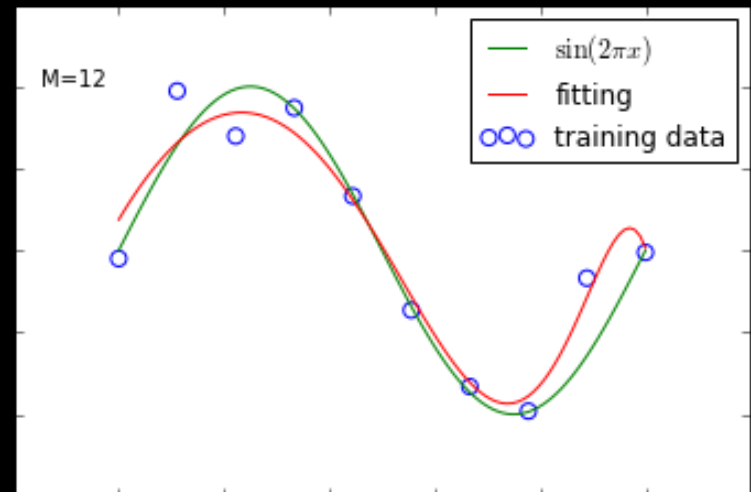
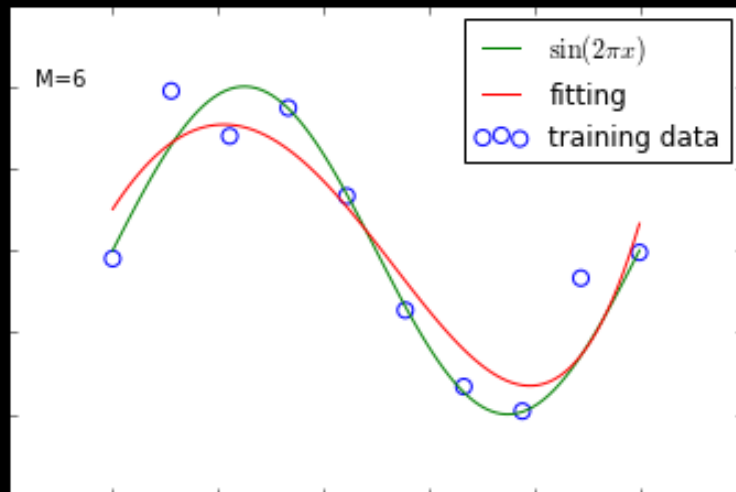
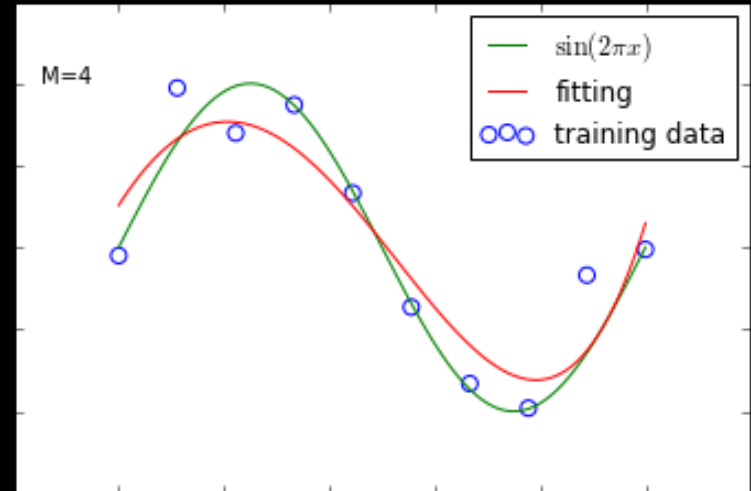
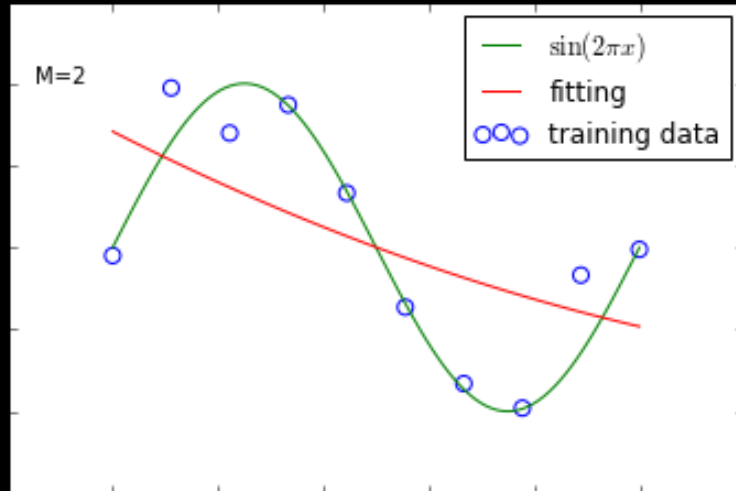
- The maximum likelihood  $\mathbf{w}$  is  $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} (\mathbf{X}^\top \mathbf{Y})$
- MAP with Gaussian prior on  $\mathbf{w}$**

$$\operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n p(y_i | \mathbf{w}, x_i) p(\mathbf{w})$$

$$y \sim \mathcal{N}(\mathbf{w} \cdot \mathbf{x}, \sigma^2) \quad \mathbf{w} \sim \mathcal{N}(0, \gamma^2)$$

$$\operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^n (y_i - \mathbf{w} x_i)^2 + \lambda \mathbf{w}^2$$

# Ridge Regression





# What we have seen

---

- **MLE with Gaussian noise is the same as minimizing the  $L_2$  error**
  - Other noise models will give other loss functions
- **MLE with a Gaussian prior adds a penalty to the  $L_2$  error, giving Ridge regression**
  - Other priors will give different penalties
- **One can make nonlinear relations linear by transforming the features**
  - Polynomial regression
  - Radial Basis Functions (RBF) – will be covered laterx

