

# Lecture 3. Linear Regression


COMP90051 Statistical Machine Learning

Lecturer: Jean Honorio



THE UNIVERSITY OF  
MELBOURNE

# This lecture

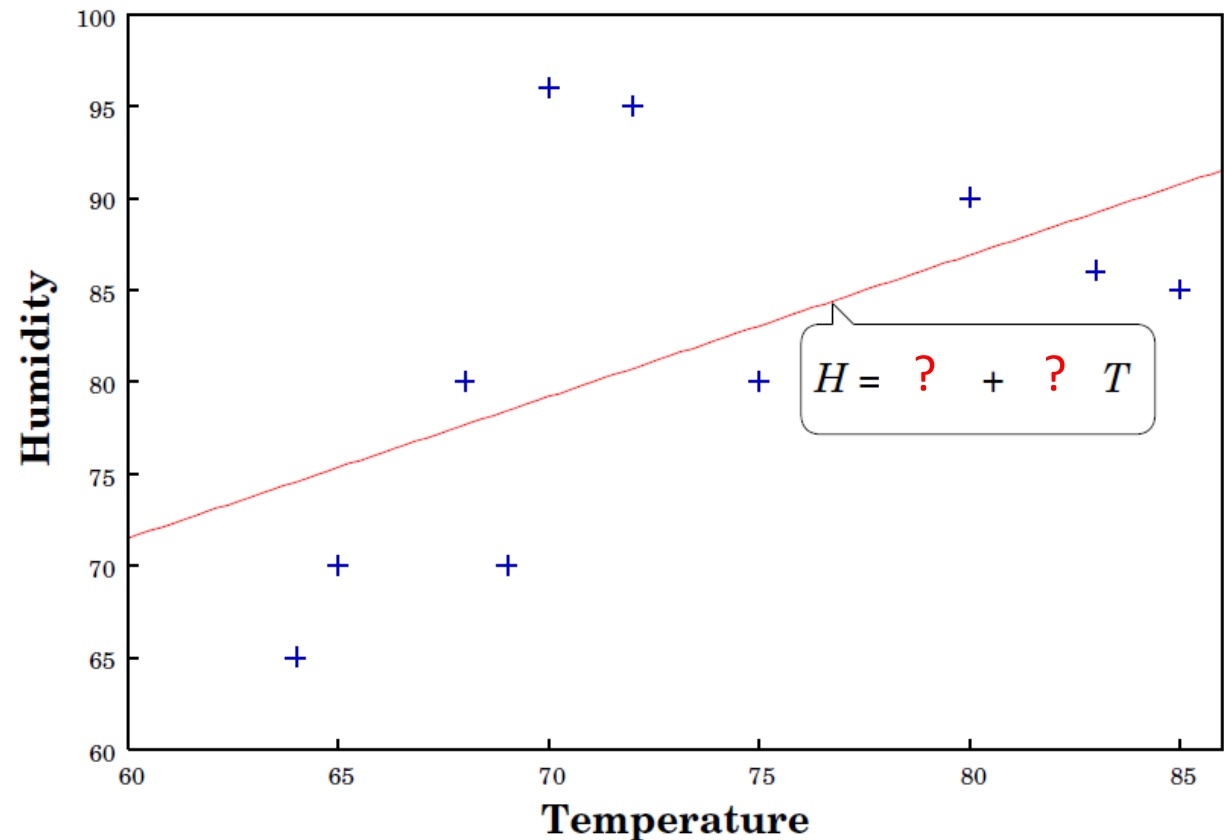
- Linear regression
  - \* Simple model (convenient maths at expense of flexibility)
  - \* Often needs less data, “interpretable”, lifts to non-linear
  - \* Derivable under all Statistical Schools: Lect 2 case study
    - This week: Frequentist + Decision theory derivations
    -  Later in semester: Bayesian approach
  - \* Convenient optimisation: Training by “analytic” (exact) solution
- Basis expansion: Data transform for more expressive models

# Linear Regression via Decision Theory

A warm-up example

# Example: Predict humidity from temperature

Temperature	Humidity
TRAINING DATA	
85	85
80	90
83	86
70	96
68	80
65	70
64	65
72	95
69	70
75	80
TEST DATA	
75	70



In regression, the task is to predict numeric response (*aka* dependent variable) from features (*aka* predictors or independent variables)

**Assume a linear relation:**  $H = a + bT$

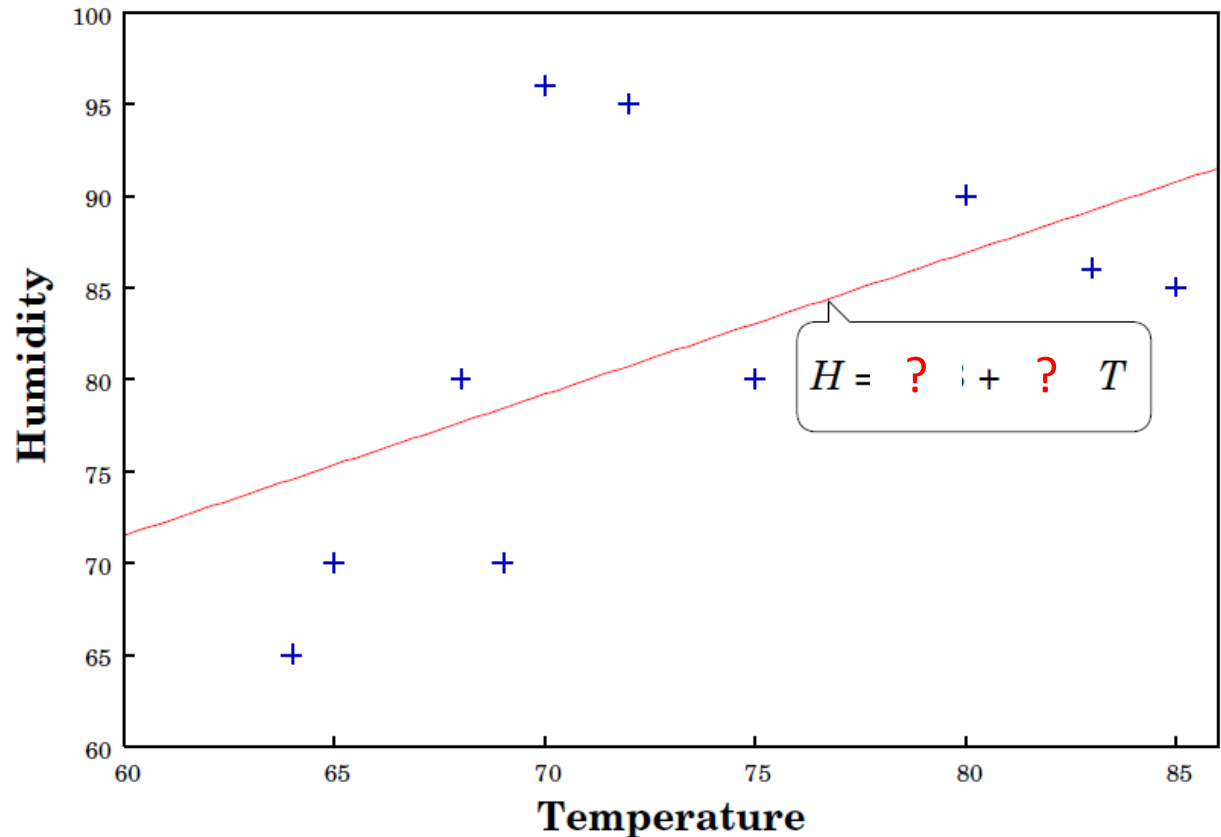
( $H$  – humidity;  $T$  – temperature;  $a, b$  – parameters)

# Example: Problem statement

- The model is  

$$H = a + bT$$
- Fitting the model = finding “best”  $a, b$  values for data at hand
- Important criterion: minimise the **sum of squared errors** (aka residual sum of squares)

*SS*  
*absolute errors*



# Example: Minimise Sum Squared Errors

To find  $a, b$  that minimise  $L = \sum_{i=1}^{10} (H_i - (a + b T_i))^2$   
 set derivatives to zero:

partial  $\frac{\partial L}{\partial a} = -2 \sum_{i=1}^{10} (H_i - a - b T_i) = 0$  objective function  
对每一个点，尽可能地拟合

if we know  $b$ , then  $\hat{a} = \frac{1}{10} \sum_{i=1}^{10} (H_i - b T_i)$  a, b 的值，找到一个最优解

$$\frac{\partial L}{\partial b} = -2 \sum_{i=1}^{10} (H_i - a - b T_i) T_i = 0$$

if we know  $a$ , then  $\hat{b} = \frac{1}{\sum_{i=1}^{10} T_i^2} \sum_{i=1}^{10} (H_i - a) T_i$

two linear equations

High-school optimisation:

- Write derivative
- Set to zero
- Solve for model
- (Check 2<sup>nd</sup> derivatives)

# Example: Analytic solution

- We have two equations and two unknowns  $a, b$
- Rewrite as a system of linear equations

$$\begin{pmatrix} 10 & \sum_{i=1}^{10} T_i \\ \sum_{i=1}^{10} T_i & \sum_{i=1}^{10} T_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^{10} H_i \\ \sum_{i=1}^{10} T_i H_i \end{pmatrix}$$

- **Analytic solution:**  $a = 25.3, b = 0.77$
- (Solve using `numpy.linalg.solve`)

→ 两个未知数两个方程

# More general decision rule

- Adopt a linear relationship between response  $y \in \mathbb{R}$  and an instance with features  $x_1, \dots, x_m \in \mathbb{R}$

$$\hat{y} = w_0 + \sum_{i=1}^m x_i w_i$$

Here  $w_0, \dots, w_m \in \mathbb{R}$  denote weights (model parameters)

- **Trick:** add a dummy feature  $x_0 = 1$  and use vector notation

$$\hat{y} = \sum_{i=0}^m x_i w_i = \mathbf{x}' \mathbf{w}$$



# Mini Summary

- Linear regression
  - \* Simple, effective, “interpretable”, basis for many approaches
  - \* Decision-theoretic frequentist derivation

Next:

Frequentist derivation; Solution/training approach

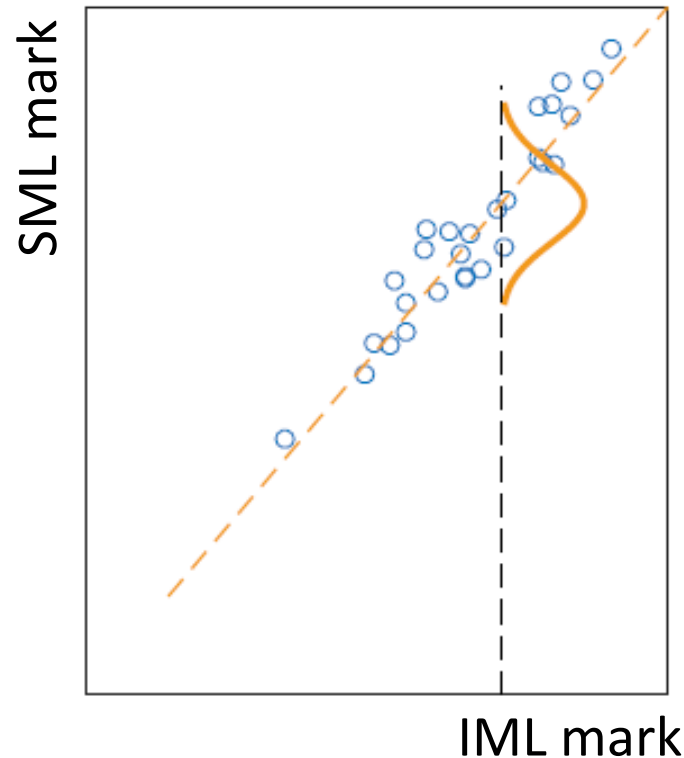
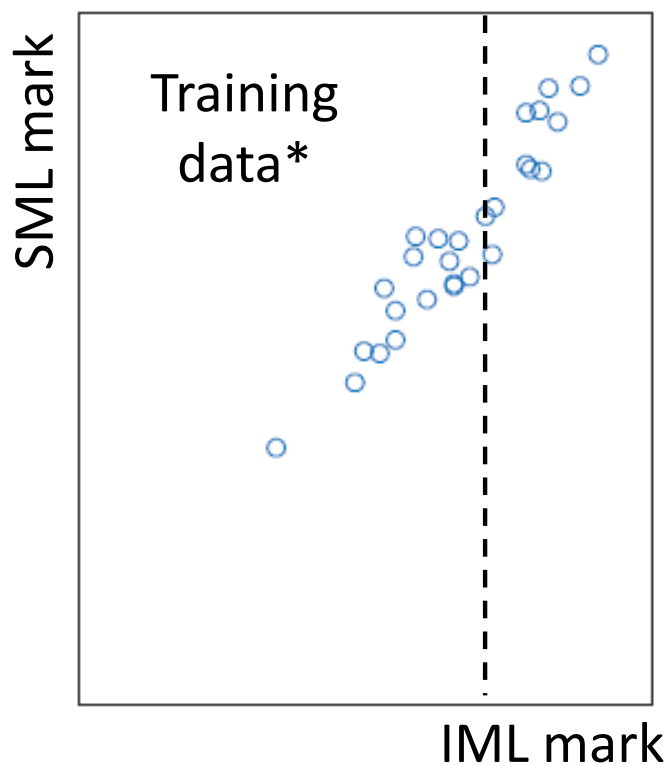
# Linear Regression via Frequentist Probabilistic Model

Max-Likelihood Estimation

# Data is noisy!

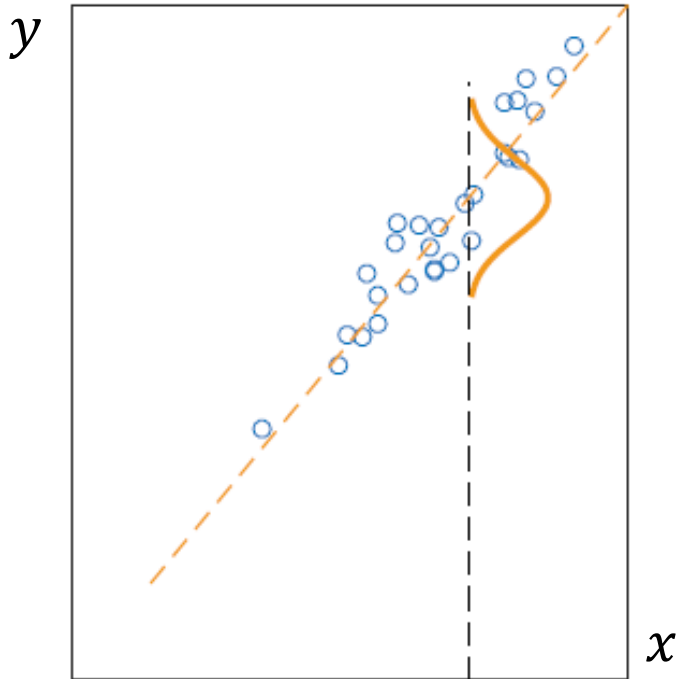
Example: predict mark for Statistical Machine Learning (SML)  
from mark for Intro ML (IML aka KT)

use input to predict output



\* synthetic data :)

# Regression as a probabilistic model



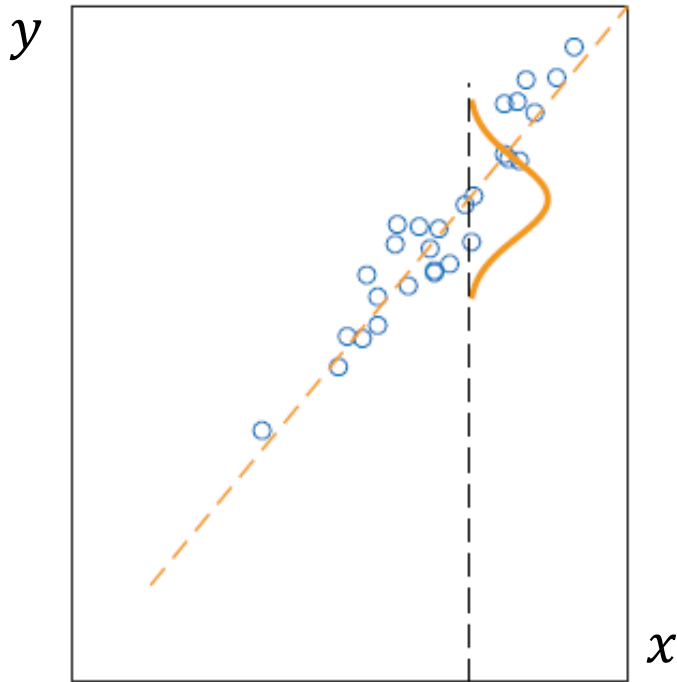
- Assume a **probabilistic model**:  
 $y = \mathbf{x}'\mathbf{w} + \varepsilon$ 
  - \* Here  $\mathbf{x}$ ,  $y$  and  $\varepsilon$  are r.v.'s
  - \* Variable  $\varepsilon$  encodes noise
- Next, assume Gaussian noise (indep. of  $\mathbf{x}$ ):  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$   
 thus:  $y \sim \mathcal{N}(\mathbf{x}'\mathbf{w}, \sigma^2)$

- Recall that  $\mathcal{N}(z; \mu, \sigma^2) \equiv \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(z-\mu)^2}{2\sigma^2}\right)$
- Therefore

$$p_{\mathbf{w}, \sigma^2}(y|\mathbf{x}) = \mathcal{N}(y; \mathbf{x}'\mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mathbf{x}'\mathbf{w})^2}{2\sigma^2}\right)$$

this is a  
squared  
error!

# Parametric probabilistic model



- Using simplified notation, **discriminative model** is: 判别式模型

$$p(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - x'w)^2}{2\sigma^2}\right)$$

条件概率


- Unknown parameters:  $w, \sigma^2$

- Given observed data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , we want to find parameter values that “best” explain the data
- Maximum-likelihood estimation**: choose parameter values that maximise the probability of observed data

# Maximum likelihood estimation

- Assuming independence of data points, the probability of data is

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n p(y_i | \mathbf{x}_i)$$

- For  $p(y_i | \mathbf{x}_i) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y_i - \mathbf{x}_i' \mathbf{w})^2}{2\sigma^2}\right)$  

- “Log trick”: Instead of maximising this quantity, we can maximise its logarithm

$$\sum_{i=1}^n \log p(y_i | \mathbf{x}_i) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i' \mathbf{w})^2 + C$$

*Handwritten Chinese notes above the equation: 求对数 (Take log), 求和 (Sum), 常数 (Constant)*

the sum of squared errors!

here  $C$  doesn't depend on  $\mathbf{w}$  (it's a constant)

- Under this model, maximising log-likelihood as a function of  $\mathbf{w}$  is equivalent to minimising the sum of squared errors

# Method of least squares

- Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ . Note bold face in  $\mathbf{x}_i$
- For convenience, place instances in rows (so attributes go in columns), representing training data as an  $n \times m$  matrix  $\mathbf{X}$ , and  $n$  vector  $\mathbf{y}$
- Probabilistic model/decision rule assumes  $\mathbf{y} \approx \mathbf{X}\mathbf{w}$
- To find  $\mathbf{w}$ , minimise the sum of squared errors

$$\begin{aligned}
 L &= \sum_{i=1}^n (y_i - \mathbf{x}_i' \mathbf{w})^2 \\
 &= \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \\
 &= (\mathbf{y} - \mathbf{X}\mathbf{w})' (\mathbf{y} - \mathbf{X}\mathbf{w}) \\
 &= \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\mathbf{w} + \underbrace{\mathbf{w}'\mathbf{X}'\mathbf{X}\mathbf{w}}_{\text{quadratic form}}
 \end{aligned}$$

$A_{2 \times 3} B_{3 \times 2} = 2 \times 2$   
 $A' B = B' A$   
 $B' A$   
 $3 \times 2 \times 2 \times 3 = 3 \times 3$   
 $\mathbf{w}' \mathbf{X}' \mathbf{y} = \mathbf{y}' \mathbf{X} \mathbf{w}$   
 since  $\|\mathbf{u}\|^2 = \mathbf{u}'\mathbf{u}$

# Method of least squares

- To find  $\mathbf{w}$ , minimise the **sum of squared errors**

$$L = \mathbf{y}'\mathbf{y} - 2\mathbf{y}'\mathbf{X}\mathbf{w} + \mathbf{w}'\mathbf{X}'\mathbf{X}\mathbf{w}$$

二次型函数:  $\frac{\partial (\mathbf{w}'\mathbf{A}\mathbf{w})}{\partial \mathbf{w}} = (\mathbf{A} + \mathbf{A}')\mathbf{w}$  (A是方阵, w是向量,  $\mathbf{A} = \mathbf{A}'$ )

$\frac{\partial (\mathbf{a}'\mathbf{X}\mathbf{w})}{\partial \mathbf{w}} = \mathbf{X}'\mathbf{a}$  (a是向量, X是矩阵, w是向量)

- Setting gradient to zero

$$\nabla L = \left[ \frac{\partial L}{\partial w_1}, \dots, \frac{\partial L}{\partial w_m} \right]' = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\mathbf{w} = \mathbf{0}$$

- Solving for  $\mathbf{w}$  yields

$$\hat{\mathbf{w}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

$(\mathbf{X}'\mathbf{X})$  is  $m \times m$ ,  $\mathbf{X}'$  is  $m \times n$ ,  $\mathbf{X}$  is  $n \times m$ ,  $\mathbf{y}$  is  $n \times 1$

Analytic solution:

- Write gradient
- Set to zero
- Solve for model

\* This system of equations called the **normal equations**

\* System is well defined only if the inverse exists



$$L = w' X' X w$$

$$\frac{\partial (w' A w)}{\partial w} = (A + A') w \quad \text{若 } A \text{ 是对称的} \Rightarrow A = A' \Rightarrow 2Aw$$

$$\frac{\partial w' (X' X) w}{\partial w} = (X' X + X' X) w = 2X' X w$$



# Bayesian derivation?

- Later in the semester: return of linear regression
- Fully Bayesian, with a posterior:
  - \* Bayesian linear regression
- Bayesian (MAP) point estimate of weight vector:
  - \* Adds a penalty term to sum of squared losses
  - \* Equivalent to  $L_2$  “regularisation” to be covered next week
  - \* Called: ridge regression

# Mini Summary

- Linear regression
    - \* Simple, effective, “interpretable”, basis for many approaches
    - \* Probabilistic frequentist derivation
    - \* Solution by normal equations
- Later in semester: Bayesian approaches

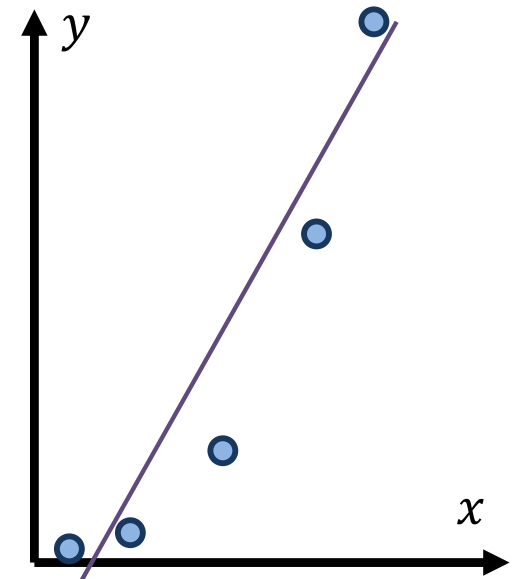
Next: Basis expansion for non-linear regression

# Basis Expansion

Extending the utility of models via  
data transformation

# Basis expansion for linear regression

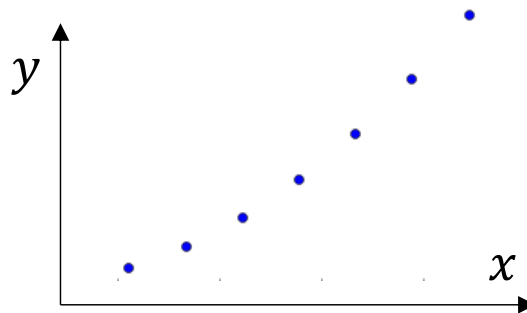
- Real data is likely to be non-linear
- What if we still wanted to use a linear regression?
  - \* Simple, easy to understand, computationally efficient, etc.
- How to marry non-linear data to a linear method?



*If you can't beat'em, join'em*

# Transform the data

- The trick is to **transform the data**: Map data into another features space, s.t. data is linear in that space
- Denote this transformation  $\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^k$ . If  $x$  is the original set of features,  $\varphi(x)$  denotes new feature set  
*transform higher dimension  $\rightarrow$  lower dimension or lower  $\rightarrow$  higher*
- Example: suppose there is just one feature  $x$ , and the data is scattered around a parabola rather than a straight line  
 *$m=1: x \Rightarrow k=2: x^2, x$*

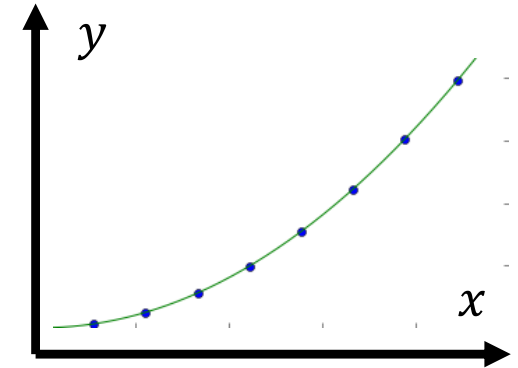


# Example: Polynomial regression

- Define

- \*  $\varphi_1(x) = x$

- \*  $\varphi_2(x) = x^2$



- Next, apply linear regression to  $\varphi_1, \varphi_2$

$$y = w_0 + w_1\varphi_1(x) + w_2\varphi_2(x) = w_0 + w_1x + w_2x^2$$

and here you have **quadratic regression**

- More generally, obtain **polynomial regression** if the new set of attributes are powers of  $x$

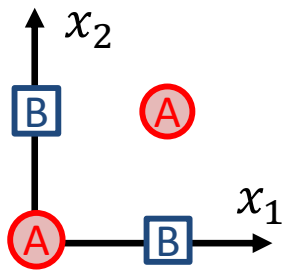
# Example: linear classification

- Example binary classification problem: Dataset not linearly separable

- Define transformation as

$$\varphi_i(\mathbf{x}) = \|\mathbf{x} - \mathbf{z}_i\|, \text{ where } \mathbf{z}_i \text{ some pre-defined constants}$$

- Choose  $\mathbf{z}_1 = [0,0]'$ ,  $\mathbf{z}_2 = [0,1]'$ ,  $\mathbf{z}_3 = [1,0]'$ ,  $\mathbf{z}_4 = [1,1]'$



there exist weights that make new data separable, e.g.:

$x_1$	$x_2$	$y$
0	0	Class A
0	1	Class B
1	0	Class B
1	1	Class A

$w_1$	$w_2$	$w_3$	$w_4$
1	0	0	1
$\varphi_1$	$\varphi_2$	$\varphi_3$	$\varphi_4$
0	1	1	$\sqrt{2}$
1	0	$\sqrt{2}$	1
1	$\sqrt{2}$	0	1
$\sqrt{2}$	1	1	0

The transformed data is linearly separable!

$\varphi'w$	$y$
$\sqrt{2}$	Class A
2	Class B
2	Class B
$\sqrt{2}$	Class A



# Radial basis functions

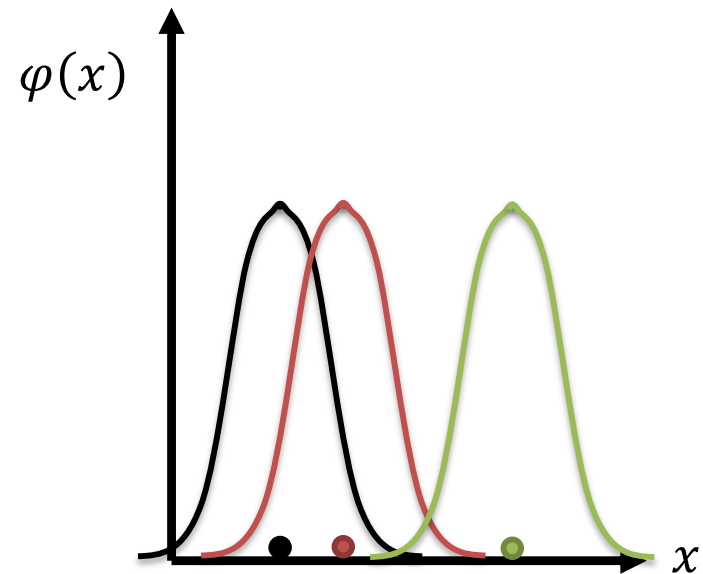
- Previous example: motivated by approximation theory where sums of RBFs approx. functions
- A **radial basis function** is a function of the form  $\varphi(\mathbf{x}) = \psi(\|\mathbf{x} - \mathbf{z}\|)$ , where  $\mathbf{z}$  is a constant

- Examples:

- $\varphi(\mathbf{x}) = \|\mathbf{x} - \mathbf{z}\|$

- $\varphi(\mathbf{x}) = \exp\left(-\frac{1}{\sigma} \|\mathbf{x} - \mathbf{z}\|^2\right)$

not linear separable  $\Rightarrow$  转换到高维可分 (separable)



# Challenges of basis expansion

- Basis expansion can significantly increase the utility of methods, especially, linear methods
- In the above examples, one limitation is that the transformation needs to be defined beforehand *rule of transformation  $z_1, z_2, z_3$* 
  - \* Need to choose the size of the new feature set *选择的个数是事先定义的*
  - \* If using RBFs, need to choose  $z_i$
- Regarding  $z_i$ , one can choose uniformly spaced points, or cluster training data and use cluster centroids
- Another popular idea is to use training data  $z_i \equiv x_i$ 
  - \* E.g.,  $\varphi_i(x) = \psi(\|x - x_i\|)$
  - \* However, for large datasets, this results in a large number of features  $\rightarrow$  computational hurdle



## Further directions

- There are several avenues for taking the idea of basis expansion to the next level
  - \* Will be covered later in this subject
- One idea is to learn the transformation  $\varphi$  from data
  - \* E.g., Artificial Neural Networks
- Another powerful extension is the use of the kernel trick
  - \* “Kernelised” methods, e.g., kernelised perceptron
- Finally, in sparse kernel machines, training depends only on a few data points
  - \* E.g., SVM

*project data on to higher dimension separable without storing the data into the computer.*