

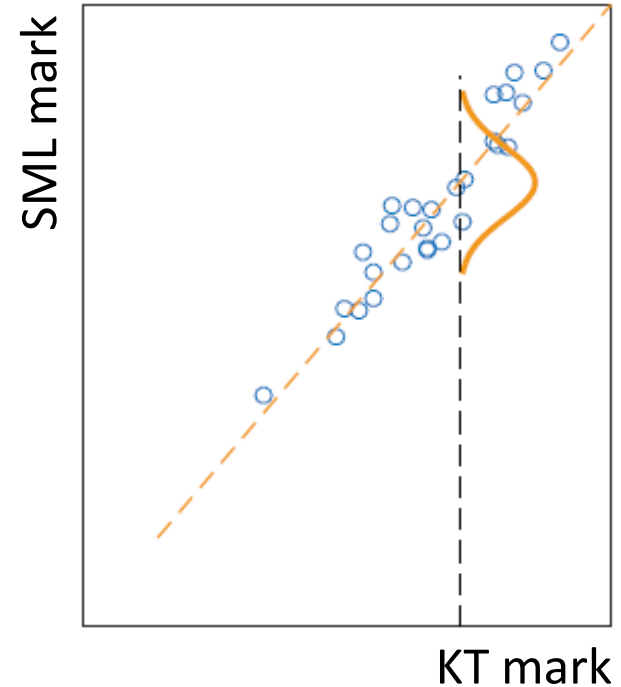
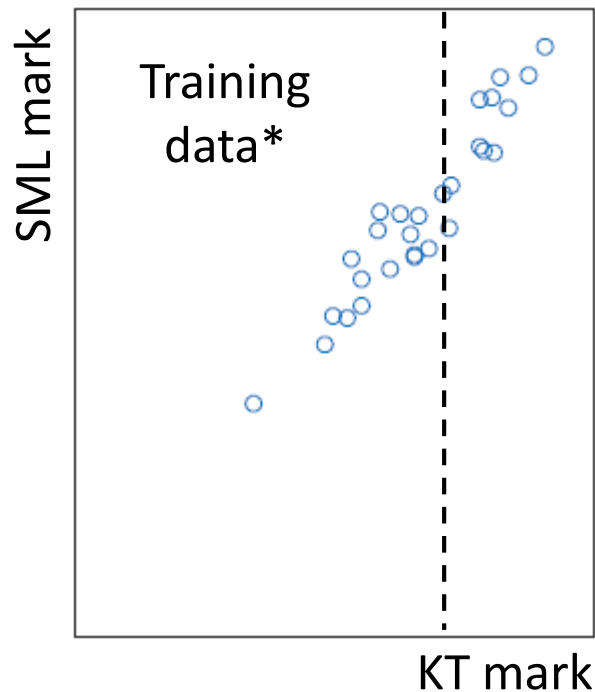


Workshop 3

COMP90051 Statistical Machine Learning
Semester 2, 2019

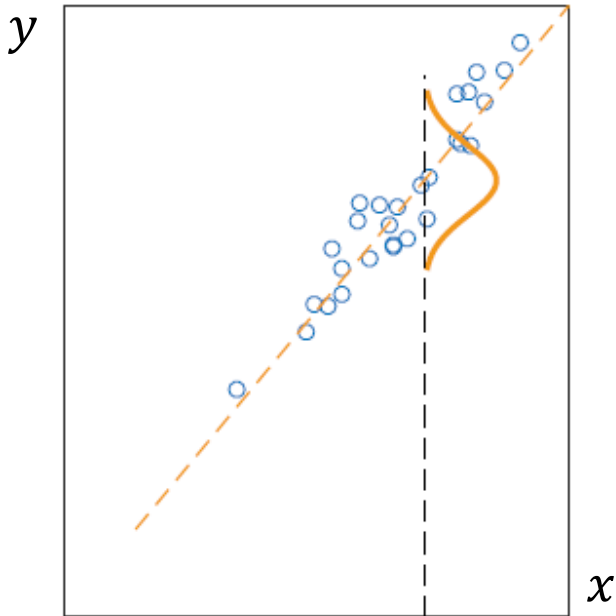
Data is noisy!

Example: predict mark for Statistical Machine Learning (SML) from mark for Knowledge Technologies (KT)



* synthetic data :)

Regression as a probabilistic model



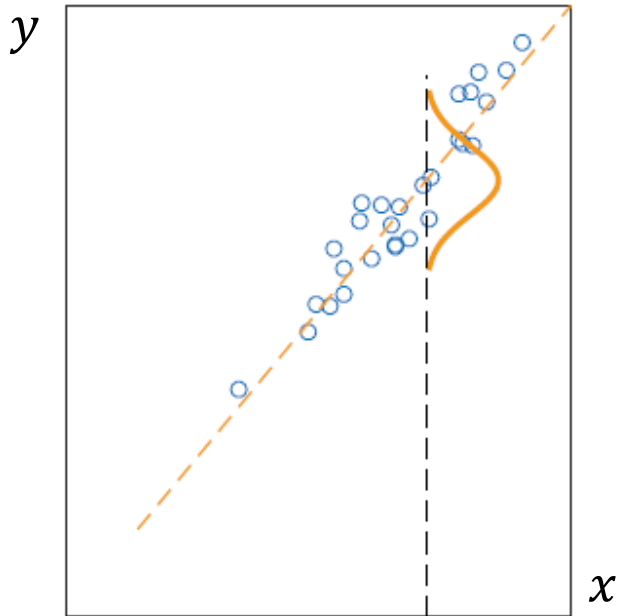
- Assume a **probabilistic model**: $Y = \mathbf{X}'\mathbf{w} + \varepsilon$
 - * Here \mathbf{X} , Y and ε are r.v.'s
 - * Variable ε encodes noise
- Next, assume Gaussian noise (indep. of \mathbf{X}):
 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

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

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

this is a squared error!

Parametric probabilistic model



- Using simplified notation, **discriminative model** is:

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

- Unknown parameters: \mathbf{w}, σ^2

- Given observed data $\{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$, 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 (why?)

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

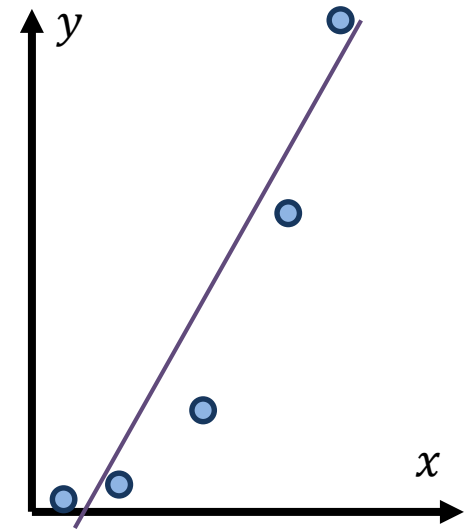
here C doesn't depend on \mathbf{w} (it's a constant)

the sum of squared errors!

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

Basis expansion for linear regression

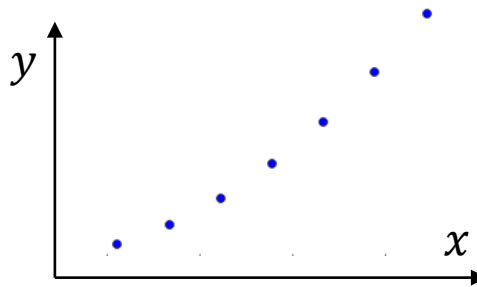
- Let's take a step back. Back to linear regression and least squares
- Real data is likely to be non-linear
- What if we still wanted to use a linear regression?
 - * It's simple, easier to understand, computationally efficient, etc.
- How to marry non-linear data to a linear method?



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

Transform the data

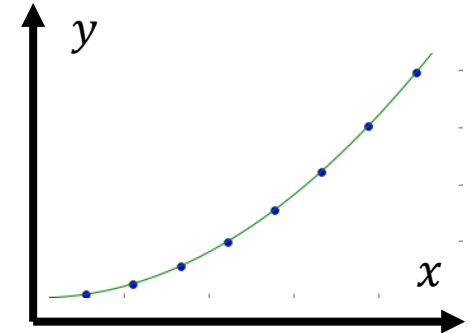
- The trick is to **transform the data**: Map data onto another features space, s.t. data is linear in that space
- Denote this transformation $\varphi: \mathbb{R}^m \rightarrow \mathbb{R}^k$. If \mathbf{x} is the original set of features, $\varphi(\mathbf{x})$ denotes new feature set
- Example: suppose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



Example: Polynomial regression

- No worries, mate: define

$$\begin{aligned}\varphi_1(x) &= x \\ \varphi_2(x) &= x^2\end{aligned}$$



- Next, apply linear regression to φ_1, φ_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

Basis expansion

- Data transformation, also known as basis expansion, is a general technique
 - * We'll see more examples throughout the course
- It can be applied for both regression and classification
- There are many possible choices of φ

