

Lecture 3b. Linear Regression - Frequentist.


COMP90051 Statistical Machine Learning

Semester 2, 2020
Lecturer: Ben Rubinstein



THE UNIVERSITY OF
MELBOURNE

This lecture

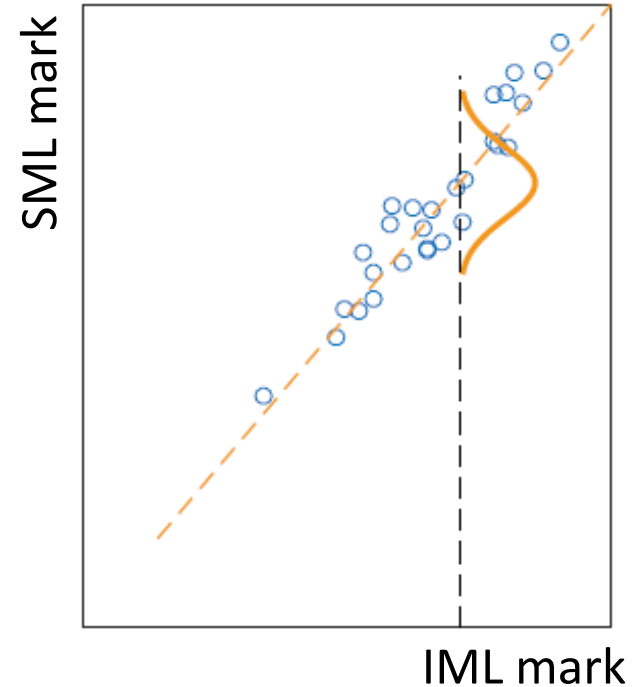
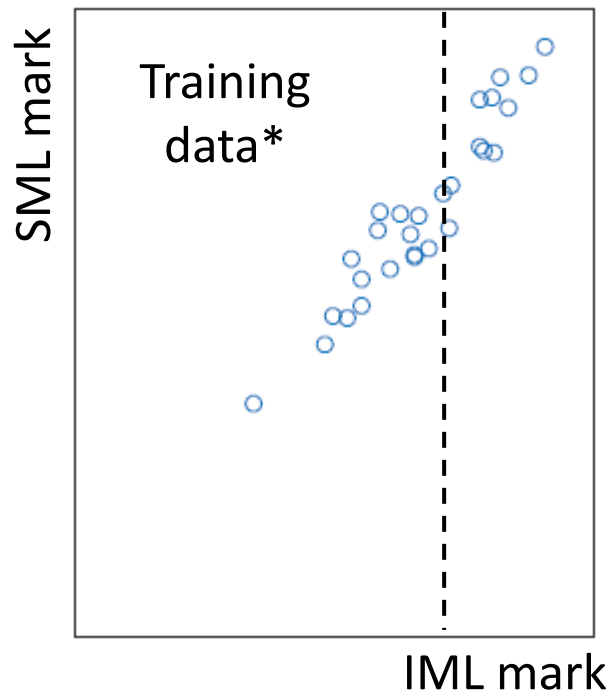
- **Linear regression**
 - * Simple model (convenient maths at expense of flexibility)
 - * Often needs less data, “interpretable”, lifts to non-linear
 - * Derivable under all Statistical Schools: Lect 2 case study
 - This week: Frequentist + Decision theory derivations
 -  Later in semester: Bayesian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

Linear Regression via Frequentist Probabilistic Model

Max-Likelihood Estimation

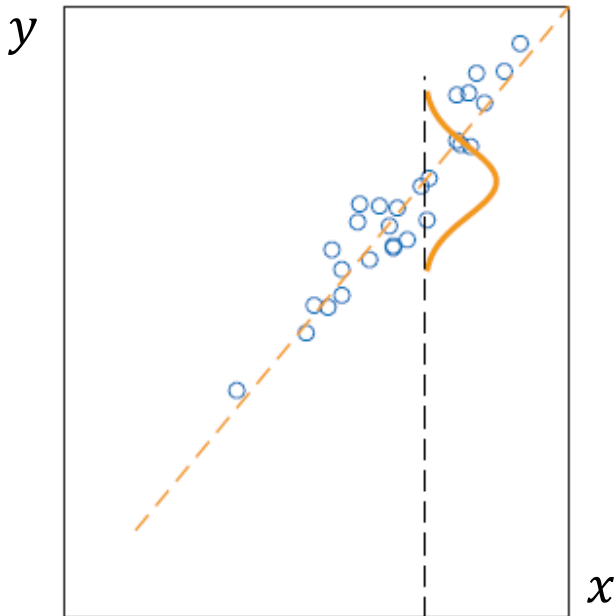
Data is noisy!

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



* synthetic data :)

Regression as a probabilistic model



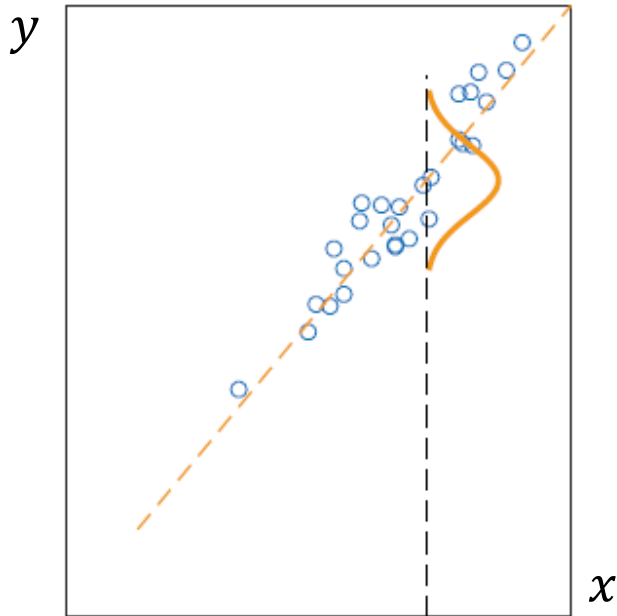
- Assume a **probabilistic model**: $Y = \mathbf{X}'\mathbf{w} + \varepsilon$
 - * Here \mathbf{X} , Y and ε are r.v.'s
 - * Variable ε encodes noise
- Next, assume Gaussian noise (indep. of \mathbf{X}):
 $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

- 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? Explained soon)

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

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

Method of least squares

Analytic solution:

- Write derivative
- Set to zero
- Solve for model

- 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 $n \times (m + 1)$ 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

$$L = \sum_{i=1}^n \left(y_i - \sum_{j=0}^m X_{ij} w_j \right)^2$$

$$L = \| \mathbf{y} - \mathbf{X}\mathbf{w} \|_2^2$$

$$\frac{\partial L}{\partial \mathbf{w}} = 2\mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

$$\mathbf{X}'\mathbf{y} = \mathbf{X}'\mathbf{X}\mathbf{w}$$

- Setting derivative to zero and solving for \mathbf{w} yields $\hat{\mathbf{w}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$

- * This system of equations called the normal equations
- * System is well defined only if the inverse exists



Wherefore art thou: 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

Summary

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

Next time: Basis expansion for non-linear regression