

Lecture 17 - Least Square Fitting

* Prayer / Quiz / AM4.

* Go over results from Exam 2

0. Numerical Calculus

* Last part of the class!

- Fitting $\hat{=}$ Interpolation \rightarrow How to turn data into functions/curves

- Integration (need interpolation to do)

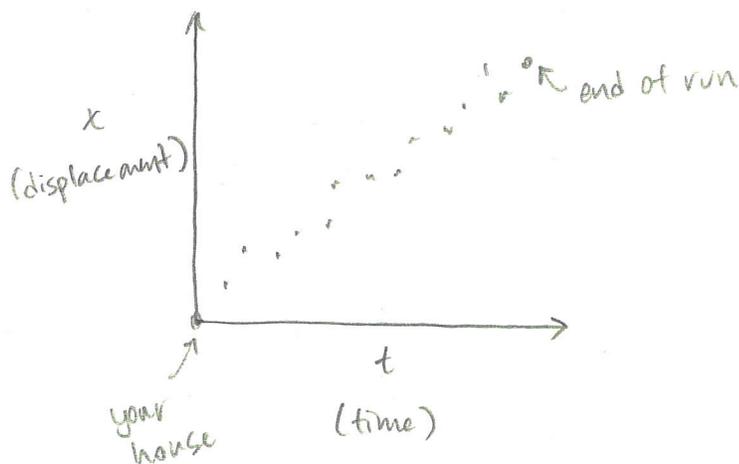
- Simple differential equations (Don't need Diff Eq. class!)

* On the home stretch. Keep doing your homework. It will payoff.

I. Least Square Fitting

* Suppose you were wearing a fitbit that tracked your position with time while you went for a run.

The data it reports might look like this



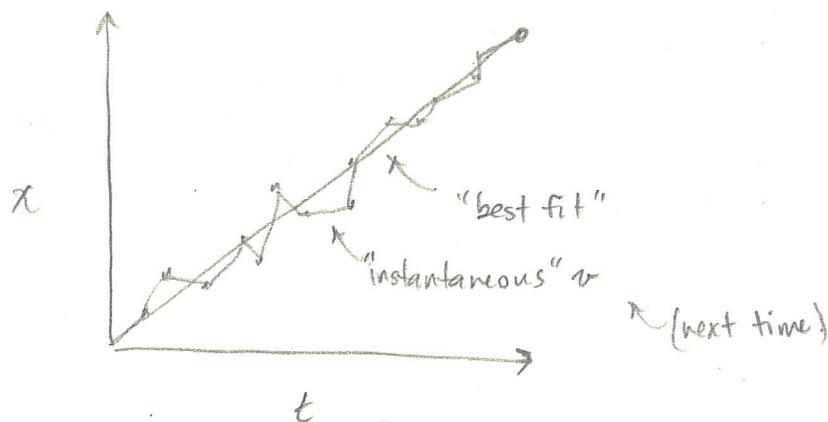
- * If we assume a mathematical model where you ran at a constant speed,

$$x = vt$$

\uparrow displacement \leftarrow time
 \leftarrow velocity

how could we find your velocity?

- * If we somehow found the "best" slope for v , that would tell us the velocity.



- * we could get little lines to connect every point, but we don't think that would tell us very much. we would have to take an average or something to find v then. Plus a lot of the spread isn't "real." It is just noise, e.g. measurement error or small variations in speed we don't care about.
- * Our model could also be wrong (here we assumed a constant velocity). Let's ignore that possibility for now

* If our model has noise, let's add that:

$$x_i = v t_i + \varepsilon_i$$

\uparrow noise ← i : data point center
 $0, 1, \dots, n-1$ points

* Least square fitting

- Vary the free parameter (in this case v) to minimize the amount of noise.

$$\varepsilon_i = x_i - v t_i$$

$$SSE = \sum_{i=0}^{n-1} (\varepsilon_i)^2 = \sum_{i=0}^{n-1} (x_i - v t_i)^2$$

SSE is a function of v . Minimize it to find v !

- This is like solving a non-linear equation with an optimization method. Here we don't expect SSE to go to zero, just that it will be as small as possible
- General formulation for $y = f(x)$:

$$\varepsilon_i = y_i - f(x_i)$$

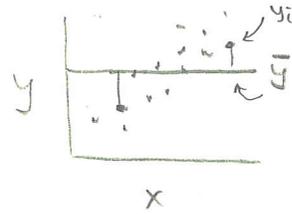
$$SSE = \sum_{i=0}^{n-1} (\varepsilon_i)^2 = \sum_{i=0}^{n-1} (y_i - f(x_i))^2$$

* Goodness of fit

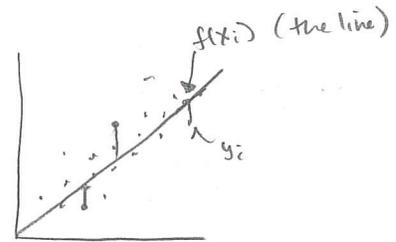
- We can tell how good the fit is by calculating the fraction of the "spread" of the data (variance) is noise or model.

Total spread: $SST = \sum_i (y_i - \bar{y})^2$

↑
average.



Spread from noise: $SSE = \sum_i (y_i - f(x_i))^2$



Fraction of spread from the model

$$R^2 = \frac{SST - SSE}{SST} \quad (\text{model} = \frac{\text{total} - \text{noise}}{\text{total}})$$

↑

called "R-squared" or the coefficient of determination.

$$R^2 = 1 - \frac{SSE}{SST}$$

II. Curve Fitting in Excel

Activity

* Demonstration & Practice of curve fitting in Excel

(1) via a trendline in a scatter plot

(2) via minimization of SSE using solver

* Calculate R^2 using RSQ

III. Curve Fitting in Python

Activity

* Demonstration & Practice of curve-fitting in Python

(1) via `numpy.polyfit()` & `numpy.polyval()`

(2) via `scipy.optimize.curve_fit()`

* calculate R^2 by hand.

If Time

Trick: Linearizing Equations

* Suppose I want to fit

$$y = A e^{Bx}$$

This isn't on my "easy list". It isn't a polynomial.

* what if I take the log of both sides:

$$\begin{aligned} \ln y &= \ln(A e^{Bx}) \\ &= \ln A + \ln e^{Bx} = \ln A + Bx \end{aligned}$$

$$\text{let } z = \ln y \Rightarrow z = A' + Bx$$

$$A' = \ln A$$

This is linear!

* Now, I can do a fit to find A' & B , then

find A :

$$A = \exp(A')$$

trend line
or
`np.polyfit()`