

# Ch 8.1: Discrete L2 and Intro to Approximation Theory

Tuesday, October 21, 2025

2:16 PM

ie - Linear Regression

aka How to be a data scientist in 5 minutes

## Context

We've seen interpolation, which has two major use cases:

- ① fit some data  $\{(x_i, y_i)\}_{i=1}^m$ , with a curve  $p$
- ② approximate a function  $f$  with a curve  $p$   
(so we choose  $x_i$ , then set  $y_i = f(x_i)$ )

Note: we index  $i, \dots, m$   
now instead of  
 $0, \dots, n$

"Interpolation" means  $p(x_i) = y_i$

Ch 8 is about Approximation Theory, with the same two scenarios!

- ① fit data. "discrete approximation" or "discrete least squares"  
or " $\ell^2$ " ("little L 2")
- ② approximate a function, "continuous approximation"  
or "continuous least squares" or " $L^2$ " ("big L 2")

...but instead of insisting  $p(x_i) = y_i$ , we only need  $p(x_i) \approx y_i$ .

Why not interpolate?

- If data are noisy, interpolation overfits to noise.
- Even if not noisy, interpolation is bad at extrapolation
- We might have physical/statistical reasons to keep  $p$  "simple"  
so it can't interpolate since not enough degrees-of-freedom.

Three types of categories:

① 1a discrete  $(x_1, x_2, \dots, x_m)$  or 1b continuous  $(\forall x \in [a, b] \dots)$

② what does  $\underbrace{p(x_i)}_{y_i} \approx y_i$  mean?

Introduce a loss function  $E$

As they come up, we'll discuss choices

③ What kind of function is  $p$ ?

e.g. polynomial  $\leftarrow$  different choices of basis

or something we can transform to a polynomial

or sines/cosines aka "trigonometric polynomials" } Fourier Series!  
any basis

Simplest scenario: discrete  $\ell^2$  "ordinary least squares" OLS

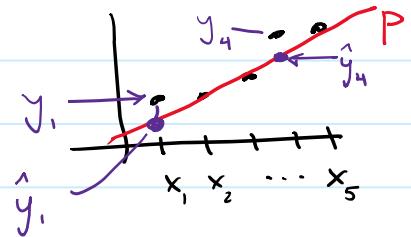
① Discrete data  $\{(x_i, y_i)\}_{i=1}^m$

② Loss is squared loss:  $E = \sum_{i=1}^m (\hat{y}_i - y_i)^2$

③ curve  $P$  is a polynomial

③a linear polynomial,  $P(x) = a_0 + a_1 x$

so  $\hat{y}_i = P(x_i)$  is our "prediction" or "estimate" of  $y_i$



So ...

$$\underset{a_0, a_1}{\text{minimize}} \quad \sum_{i=1}^m \underbrace{(a_0 + a_1 x_i - y_i)^2}_{\hat{y}_i} \quad \} E$$

We can solve w/ linear algebra! Minimizers must satisfy

$\frac{\partial E}{\partial a_0} = 0$  and  $\frac{\partial E}{\partial a_1} = 0$ . In fact, since  $E$  is jointly convex in  $a_0$  and  $a_1$ , these conditions are also sufficient.

$$(*) \quad \frac{\partial E}{\partial a_0} = \sum_{i=1}^m 2 \cdot (a_0 + a_1 x_i - y_i) = 0$$

$$(*) \quad \frac{\partial E}{\partial a_1} = \sum_{i=1}^m 2 x_i (a_0 + a_1 x_i - y_i) = 0$$

2 equations, linear in  $a_0$  and  $a_1$ .

"normal"  
here refers  
to "orthogonal"

Rewrite in vector form:

$$\vec{y} = (y_i) \in \mathbb{R}^m, \quad X = \begin{bmatrix} x_0 & x_1 \\ \vdots & \vdots \\ x_m & x_1 \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_m \end{bmatrix} \in \mathbb{R}^{m \times 2}, \quad \vec{a} = [a_0 \ a_1] \in \mathbb{R}^2$$

$$\text{so } \hat{y} = X \cdot \vec{a}$$

$$\text{and we want } \min_{\vec{a} \in \mathbb{R}^2} \|X \cdot \vec{a} - \vec{y}\|_2^2 \quad \text{since } \|r\|_2 = \sqrt{\sum_i r_i^2}$$

$$= \underbrace{(X \vec{a} - \vec{y})^T (X \vec{a} - \vec{y})}_{= \vec{a}^T X^T X \vec{a} - 2 \vec{y}^T X \vec{a} + \vec{y}^T \vec{y}}$$

Solution  $\vec{a}$  satisfies the normal equations

$$\underbrace{X^T X}_{2 \times 2} \underbrace{\vec{a}}_{2 \times 1} = \underbrace{X^T \vec{y}}_{2 \times 1}$$

equivalent to (\*)

Tattoo this to your arm.  
Don't forget it.

Variations:

① non-linear polynomials, eg.  $p(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3$  say, degree  $n$

then set  $X = \begin{bmatrix} x_0 & x_1 & x_2 & x_3 \\ \vdots & \vdots & \vdots & \vdots \\ x_m^0 & x_m^1 & x_m^2 & x_m^3 \end{bmatrix} \in \mathbb{R}^{m \times n+1}$ ,  $\vec{a} = \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{bmatrix} \in \mathbb{R}^{n+1}$

and solve  $\underline{\underline{X}}^T X \vec{a} = \underline{\underline{X}}^T \vec{y}$  as before

Comments:

① Does  $\underline{\underline{X}}^T X \vec{a} = \underline{\underline{X}}^T \vec{y}$  have a solution? Is it unique?  
Gram matrix

If points  $x_i$  are distinct and  $n+1 \leq m$  then  
columns of  $X$  are linearly independent, so  $\underline{\underline{X}}^T X$  is  
positive definite, hence there is a solution and it's unique  
( $\Rightarrow$  invertible)

②  $X$  is a Vandermonde type matrix and ill-conditioned  
if  $n$  is large.

Two tricks: ① Choose a better polynomial basis.  
More on this later this chapter

② Don't form  $\underline{\underline{X}}^T X$  since that exacerbates  
the ill-conditioning. More on this in ch. 6  
(e.g., QR or SVD) Short answer: use  
np.linalg.lstsq( $X, y$ )

② change-of-variables

ex.  $y = a_0 \cdot b^{a_1 x}$  or  $y = a_0 x^{a_1}$

take  $\log(\cdot)$  of these models: make it linear

rename  $\underline{\underline{y}}$   $\underline{\underline{a}_0}$   $\underline{\underline{a}_1}$   $\underline{\underline{x}}$  or  $\underline{\underline{y}}$   $\underline{\underline{a}_0}$   $\underline{\underline{a}_1}$   $\underline{\underline{x}}$

Solve for new variables,

then transform back

⚠ The transformation affects how we interpret our loss

OR ... keep it nonlinear and solve non-linear least squares using ch. 10 techniques

## Other tricks

model:  $y = \alpha_1 \cdot \sin(\omega t + \alpha_2)$

Known frequency  
fit amplitude ...and phase

Trick  $\sin(u+v) = \underline{\sin(u)} \cdot \cos(v) + \underline{\cos(u)} \cdot \sin(v)$

so fit  $y = \tilde{\alpha}_1 \cdot \underline{\sin(\omega t)} + \tilde{\alpha}_2 \cdot \underline{\cos(\omega t)}$  linear in parameters!

$$\left. \begin{array}{l} \tilde{\alpha}_1 = \alpha_1 \cdot \cos(\alpha_2) \\ \tilde{\alpha}_2 = \alpha_1 \cdot \sin(\alpha_2) \end{array} \right\}$$

solve:  $\tilde{\alpha}_1^2 + \tilde{\alpha}_2^2 = \alpha_1^2 (\cos^2 + \sin^2) = \alpha_1^2$

to find  $\alpha_1$ ,

then  $\alpha_2 = \cos^{-1}(\tilde{\alpha}_1/\alpha_1)$ .

nonlinear: `scipy.optimize.least_squares`

linear: `scipy.linalg.lstsq` Don't use if you don't need to!

### (3) Weighted & Generalized Least Squares

if  $E = \sum_{i=1}^m w_i (\hat{y}_i - y_i)^2$  (eg.  $w_i$  reflects an estimate of accuracy of  $i^{th}$  data point)

then if  $W = \text{diag}(\vec{w})$ ,

$$\text{solve } \min_{\vec{\alpha}} (X\vec{\alpha} - \vec{y})^T W (X\vec{\alpha} - \vec{y})$$

$$w_1 \text{ solution } X^T W X \vec{\alpha} = X^T W \vec{y} \quad \text{Weighted least squares}$$

even more generally, let  $C > 0$  be positive definite,

$$\text{solve } \min_{\vec{\alpha}} (X\vec{\alpha} - \vec{y})^T C (X\vec{\alpha} - \vec{y})$$

$$\underbrace{\|X\vec{\alpha} - \vec{y}\|_C^2}_{\text{"Mahalanobis distance"}}$$

w<sub>1</sub> solution

$$X^T C X \vec{\alpha} = X^T C \vec{y} \quad \text{Generalized Least Squares}$$

Ex:  $C = \text{estimate of inverse of covariance matrix}$

⚠ Careful w/ terminology.

"General linear model" is for multiple RHS  $\vec{y}$

"Generalized linear model" if you have a "link" function, as in logistic regression.

Neither is the same as Generalized L.S.

(4) (4a)  $E = \sum_{i=1}^m |\hat{y}_i - y_i|$  "least absolute deviation" aka  $\min_{\vec{\alpha}} \|X\vec{\alpha} - \vec{y}\|_1$ ,

(4b)  $E = \max_{1 \leq i \leq m} |\hat{y}_i - y_i|$  "minimax" aka  $\min_{\vec{\alpha}} \|X\vec{\alpha} - \vec{y}\|_\infty$

Both reasonable, depend on where you think the "error" in data comes from. We'll explore in lab. Both are a special type of convex optimization problem known as a linear program. We've known how to solve since 1940's ... but Gauss, etc., couldn't solve. It's not just linear algebra.

One way these  $\lambda$ , and  $\lambda \infty$  error terms arise is in **maximum likelihood estimation (MLE)**, a very common statistical technique.

MODEL:

$$y = \hat{y}_\theta + \text{noise}$$

↑  
we'll observe      ↑  
our model,  
parameterized by  $\theta$

↗ unobserved but we might know its  
distribution

Ex:  $\hat{y}_\theta = \theta_0 + \theta_1 x + \theta_2 x^2$

Ex: assume noise  $\sim N(0, \sigma^2)$  "normal" / "Gaussian"

so the probability density function (pdf)

$$\text{is } P_\theta(y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\hat{y}_\theta)^2/2\sigma^2}$$

abbreviate ↗  
lecture

If we have multiple observations  $\{y_i\}$  that are independent,

$$\text{then } P_\theta(\{y_i\}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-{(y_i - \hat{y}_{i,\theta})^2}/2\sigma^2)$$

The "likelihood" is just  $P_\theta(\{y_i\})$  but as a function of  $\theta$

Given a value of  $\theta$ , it is "how likely would it be to see these observations if  $\theta$  were the model parameters"

So, find the  $\theta$  to maximize the likelihood... hence MLE.

Note:  $P_\theta(\dots) \geq 0$  so we could take  $\log(P_\theta(\dots))$  w/o changing location of maximizer.      monotonic

$$\text{So, } \max_{\theta} P_\theta(\{y_i\}) \underset{\theta}{\approx} \min_{\theta} -\log(P_\theta(\{y_i\}))$$

negative log-likelihood

if noise  $\stackrel{iid}{\sim} N(0, \sigma^2)$ ,

$$\text{this is } \min_{\theta} -\log \left( \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-{(y_i - \hat{y}_{i,\theta})^2}/2\sigma^2) \right)$$

$$\text{i.e. } \min_{\theta} \sum_{i=1}^m (y_i - \hat{y}_{i,\theta})^2 + \text{constants}$$

no  $\theta$  dependence

i.e. least squares!

↗ So popular because:

(1) it's MLE for Gaussian noise, and Gaussians are common (ex: CLT)

(2) it's easy to solve and well understood

## ⑤ Regularized least squares ... in particular, Tikhonov / Ridge Regression

$$\min_{\vec{\alpha}} \|\mathbf{X}\vec{\alpha} - \vec{y}\|_2^2 + \lambda \cdot \|\vec{\alpha}\|_2^2$$

$\lambda > 0$  regularization term  
a hyperparameter

Still has closed form solution:

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \vec{\alpha} = \mathbf{X}^T \vec{y}$$

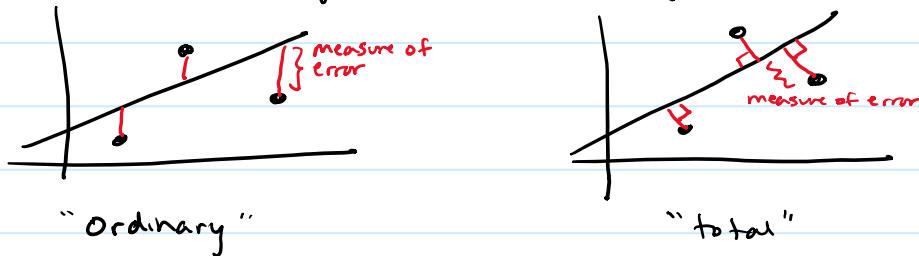
Many variants beyond  $\|\vec{\alpha}\|_2^2$

Arises in Maximum a posteriori (MAP) statistical estimation, which is a (kind of) Bayesian technique that assumes a "prior" distribution on  $\vec{\alpha}$  (not "fully" Bayesian since it's just a point estimate)

Tikhonov corresponds to assuming  $\vec{\alpha}$  is Gaussian.

## ⑥ Total Least Squares + Orthogonal distance regression

Ordinary regression assumes we have error in our response variable  $y$   
... but what if you have error in your independent variable  $x$ ?



Can solve w/ SVD. We won't go into detail

Scipy.odr solves. You should provide it w/ an estimate of the variance in  $y$  and variance in  $x$

$\underbrace{\text{if exactly equal, then error is exactly the perpendicular}}$