

18-661 Introduction to Machine Learning

Linear Regression – Part I

Spring 2023

ECE – Carnegie Mellon University

Linear Regression

Recap of MLE/MAP

Linear Algebra Review

Linear Regression

- Formulation

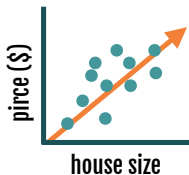
- Univariate Solution

- Multivariate Solution

- Probabilistic Interpretation

Task 1: Regression

How much should you sell your house for?



input: houses & features **learn:** $x \rightarrow y$ relationship **predict:** y (*continuous*)

Course Covers: Linear/Ridge Regression, Loss Function, SGD, Feature Scaling, Regularization, Cross Validation

Supervised learning

In a supervised learning problem, you have access to input variables (X) and outputs (Y), and the goal is to predict an output given an input

- Examples:
 - **Housing prices (Regression)**: predict the price of a house based on features (size, location, etc)
 - **Cat vs. Dog (Classification)**: predict whether a picture is of a cat or a dog

Predicting a continuous outcome variable:

- Predicting a company's future stock price using its profit and other financial info
- Predicting annual rainfall based on local flora and fauna
- Predicting distance from a traffic light using LIDAR measurements

Magnitude of the error matters:

- We can measure 'closeness' of prediction and labels, leading to different ways to evaluate prediction errors.
 - Predicting stock price: better to be off by 1\$ than by 20\$
 - Predicting distance from a traffic light: better to be off 1 m than by 10 m
- We should choose learning models and algorithms accordingly.

Predicting House Prices: Collecting Data

[Overview](#)
[Property Details](#)
[Tour Insights](#)
[Property History](#)
[Public Records](#)
[Activity](#)
[Schools](#)

Five unit apartment complex within 2 blocks of USC campus, Gate #6. Great for students (most student leases have parents as guarantors). Most USC students live off campus, so housing units like this are always fully leased. Situated on a gated, corner lot, and across from an elementary school, this complex was recently renovated, and has in-unit laundry hook ups, wall-unit AC, and 1/2 parking spaces. It's within a (DPS Department of Public Safety) and Campus Cruiser patrolled area. This is a great income generating property, not to be missed!

Property Type: Multi-Family
Community: Downtown Los Angeles
MLS#: 22176741

Style: Two Level, Low Rise
County: [Los Angeles](#)

Property Details for 3620 South BUDLONG, Los Angeles, CA 90007

Details provided by iTech MLS and may not match the public record. [Learn More](#)

Interior Features

Kitchen Information

- Remodeled
- Oven, Range

Laundry Information

- Inside Laundry

Heating & Cooling

- Inside Cooling Unit(s)

Multi-Unit Information

Community Features

- Units In Complex (Total): 5

Multi-Family Information

- # Leased: 5
- # of Buildings: 1
- Owner Pays Water
- Tenant Pays Electricity, Tenant Pays Gas

Unit 1 Information

- # of Beds: 2
- # of Baths: 1
- Unfurnished
- Monthly Rent: \$1,700

Unit 2 Information

- # of Beds: 3
- # of Baths: 1
- Unfurnished
- Monthly Rent: \$2,250

Unit 3 Information

- Unfurnished

Unit 4 Information

- # of Beds: 3
- # of Baths: 1
- Unfurnished

Monthly Rent: \$2,350

Unit 5 Information

- # of Beds: 3
- # of Baths: 2
- Unfurnished
- Monthly Rent: \$2,325

Unit 6 Information

- # of Beds: 3
- # of Baths: 1
- Monthly Rent: \$2,250

Property / Lot Details

Property Features

- Automatic Gate, Card/Code Access

- Automatic Gate, Lawn, Sidewalks
- Corner Lot, Near Public Transit

- Tax Parcel Number: 5042017019

Lot Information

- Lot Size (Sq. Ft.): 9,648
- Lot Size (Acres): 0.2215
- Lot Size Source: Public Records

Property Information

- Updated/Remodeled
- Square Footage Source: Public Records

Parking / Garage, Exterior Features, Utilities & Financing

Parking Information

- # of Parking Spaces (Total): 12
- Parking Space
- Gated

Utility Information

- Green Certification Rating: 0.00
- Green Location: Transportation, Walkability
- Green Walk Score: 0
- Green Year Certified: 0

Financial Information

- Capitalization Rate (%): 8.25
- Actual Annual Gross Rent: \$128,331
- Gross Rent Multiplier: 11.29

Location Details, Misc. Information & Listing Information

Location Information

- Cross Street: W 36th Pl

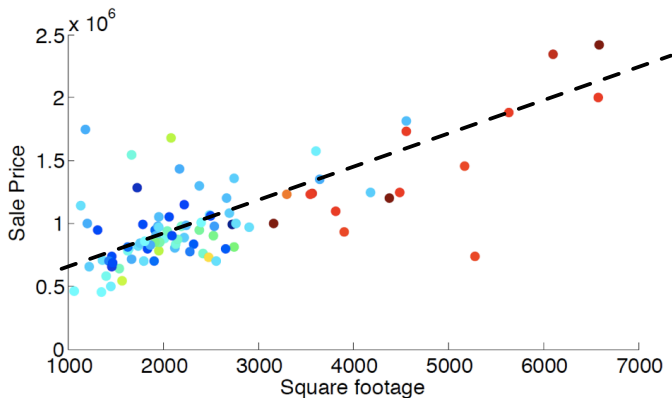
Expense Information

- Operating: \$37,984

Listing Information

- Listing Terms: Cash, Cash To Existing Loan
- Buyer Financing: Cash

Correlation between Square Footage and Sale Price



- Sale price \approx price_per_sqft \times square_footage + fixed_expense
- Learn parameters (w_0 , w_1) of the dotted line $y = w_1x + w_0$

Data from Many Houses?

- $\text{Sale_price} = \text{price_per_sqft} \times \text{square_footage} + \text{fixed_expense} + \text{unexplainable_stuff}$
- Want to learn the price_per_sqft and fixed_expense
- **Training data:** past sales records

| sqft | sale price |
|------|------------|
| 2000 | 800K |
| 2100 | 907K |
| 1100 | 312K |
| 5500 | 2,600K |
| ... | ... |

Problem: there isn't a $\mathbf{w} = [w_1, w_0]^T$ that will satisfy all equations

Reduce Prediction Error

How to measure errors?

| sqft | sale price | prediction | abs error | squared error |
|------|------------|------------|-----------|---------------|
| 2000 | 810K | 720K | 90K | 8100 |
| 2100 | 907K | 800K | 107K | 107^2 |
| 1100 | 312K | 350K | 38K | 38^2 |
| 5500 | 2,600K | 2,600K | 0 | 0 |
| ... | ... | | | |

- **absolute** difference (ℓ_1 norm): $|\text{prediction} - \text{sale price}|$.
- **squared** difference (ℓ_2 norm): $(\text{prediction} - \text{sale price})^2$
[differentiable!].

Minimize Squared Errors

Our model:

Sale_price =

price_per_sqft \times square_footage + fixed_expense + unexplainable_stuff

Training data:

| sqft | sale price | prediction | error | squared error |
|-------|------------|------------|-------|-----------------------------------|
| 2000 | 810K | 720K | 90K | 8100 |
| 2100 | 907K | 800K | 107K | 107^2 |
| 1100 | 312K | 350K | 38K | 38^2 |
| 5500 | 2,600K | 2,600K | 0 | 0 |
| ... | ... | | | |
| Total | | | | $8100 + 107^2 + 38^2 + 0 + \dots$ |

Aim:

Adjust price_per_sqft and fixed_expense such that the sum of the squared error is minimized — i.e., the unexplainable_stuff is minimized.

Linear Regression

Setup:

- **Input:** $\mathbf{x} \in \mathbb{R}^D$ (covariates, predictors, features, etc)
- **Output:** $y \in \mathbb{R}$ (responses, targets, outcomes, outputs, etc)
- **Model:** $f: \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_{d=1}^D w_d x_d = w_0 + \mathbf{w}^\top \mathbf{x}$.
 - $\mathbf{w} = [w_1 \ w_2 \ \cdots \ w_D]^\top$: *weights, parameters, or parameter vector*
 - w_0 is called *bias*.
 - Sometimes, we also call $\tilde{\mathbf{w}} = [w_0 \ w_1 \ w_2 \ \cdots \ w_D]^\top$ parameters.
- **Training data:** $\mathcal{D} = \{(\mathbf{x}_n, y_n), n = 1, 2, \dots, N\}$

Minimize the Residual Sum of Squares:

$$RSS(\tilde{\mathbf{w}}) = \sum_{n=1}^N [y_n - f(\mathbf{x}_n)]^2 = \sum_{n=1}^N [y_n - (w_0 + \sum_{d=1}^D w_d x_{nd})]^2$$

Recap of MLE/MAP

Linear Algebra Review

Linear Regression

Formulation

Univariate Solution

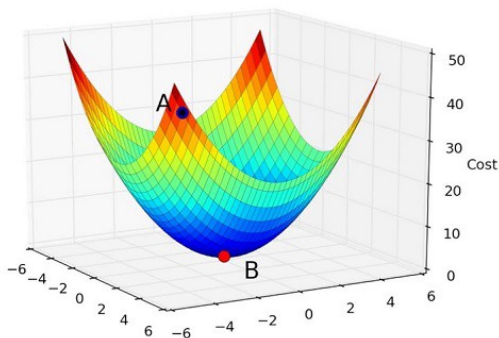
Multivariate Solution

Probabilistic Interpretation

A Simple Case: x Is One-dimensional ($D=1$)

Residual sum of squares:

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + w_1 x_n)]^2$$



What kind of function is this? CONVEX (has a unique global minimum)

A Simple Case: x Is One-dimensional ($D=1$)

Residual sum of squares:

$$RSS(\mathbf{w}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + w_1 x_n)]^2$$

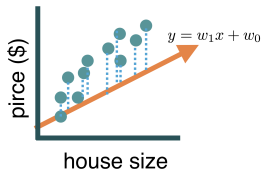


Figure 2: RSS is the sum of squares of the dotted lines

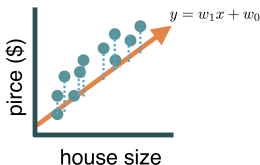


Figure 3: Adjust (w_0, w_1) to reduce RSS

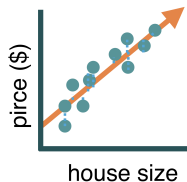


Figure 4: RSS minimized at (w_0^*, w_1^*)

A Simple Case: x Is One-dimensional ($D=1$)

Residual sum of squares:

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + w_1 x_n)]^2$$

Stationary points:

Take derivative with respect to parameters and set it to zero

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_0} = 0 \Rightarrow -2 \sum_n [y_n - (w_0 + w_1 x_n)] = 0,$$

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_1} = 0 \Rightarrow -2 \sum_n [y_n - (w_0 + w_1 x_n)] x_n = 0.$$

A Simple Case: x Is One-dimensional ($D=1$)

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_0} = 0 \Rightarrow -2 \sum_n [y_n - (w_0 + w_1 x_n)] = 0$$

$$\frac{\partial RSS(\tilde{\mathbf{w}})}{\partial w_1} = 0 \Rightarrow -2 \sum_n [y_n - (w_0 + w_1 x_n)] x_n = 0$$

Simplify these expressions to get the “Normal Equations”:

$$\sum y_n = N w_0 + w_1 \sum x_n$$

$$\sum x_n y_n = w_0 \sum x_n + w_1 \sum x_n^2$$

Solving the system we obtain the **least squares coefficient estimates**:

$$w_1 = \frac{\sum (x_n - \bar{x})(y_n - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{and} \quad w_0 = \bar{y} - w_1 \bar{x}$$

where $\bar{x} = \frac{1}{N} \sum_n x_n$ and $\bar{y} = \frac{1}{N} \sum_n y_n$.

Example

| sqft (1000's) | sale price (100k) |
|---------------|-------------------|
| 1 | 2 |
| 2 | 3.5 |
| 1.5 | 3 |
| 2.5 | 4.5 |

Residual sum of squares:

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - f(\mathbf{x}_n)]^2 = \sum_n [y_n - (w_0 + w_1 x_n)]^2$$

The w_1 and w_0 that minimize this are given by:

$$w_1 \approx 1.6$$

$$w_0 \approx 0.45$$

Recap of MLE/MAP

Linear Algebra Review

Linear Regression

Formulation

Univariate Solution

Multivariate Solution

Probabilistic Interpretation

Least Mean Squares: \mathbf{x} Is D -dimensional

| sqft (1000's) | bedrooms | bathrooms | sale price (100k) |
|---------------|----------|-----------|-------------------|
| 1 | 2 | 1 | 2 |
| 2 | 2 | 2 | 3.5 |
| 1.5 | 3 | 2 | 3 |
| 2.5 | 4 | 2.5 | 4.5 |

RSS($\tilde{\mathbf{w}}$) in matrix form:

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_n]^2,$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^\top, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^\top$$

What is $\tilde{\mathbf{x}}$ for the first house? $[1, 1, 2, 1]^\top$

Least Mean Squares: \mathbf{x} Is D -dimensional

$RSS(\tilde{\mathbf{w}})$ in matrix form:

$$RSS(\tilde{\mathbf{w}}) = \sum_n [y_n - (w_0 + \sum_d w_d x_{nd})]^2 = \sum_n [y_n - \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_n]^2,$$

where we have redefined some variables (by augmenting)

$$\tilde{\mathbf{x}} \leftarrow [1 \ x_1 \ x_2 \ \dots \ x_D]^\top, \quad \tilde{\mathbf{w}} \leftarrow [w_0 \ w_1 \ w_2 \ \dots \ w_D]^\top$$

which leads to

$$\begin{aligned} RSS(\tilde{\mathbf{w}}) &= \sum_n (y_n - \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_n)(y_n - \tilde{\mathbf{x}}_n^\top \tilde{\mathbf{w}}) \\ &= \sum_n \tilde{\mathbf{w}}^\top \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^\top \tilde{\mathbf{w}} - 2y_n \tilde{\mathbf{x}}_n^\top \tilde{\mathbf{w}} + \text{const.} \\ &= \left\{ \tilde{\mathbf{w}}^\top \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^\top \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^\top \right) \tilde{\mathbf{w}} \right\} + \text{const.} \end{aligned}$$

RSS($\tilde{\mathbf{w}}$) in New Notations

From previous slide:

$$RSS(\tilde{\mathbf{w}}) = \left\{ \tilde{\mathbf{w}}^\top \left(\sum_n \tilde{\mathbf{x}}_n \tilde{\mathbf{x}}_n^\top \right) \tilde{\mathbf{w}} - 2 \left(\sum_n y_n \tilde{\mathbf{x}}_n^\top \right) \tilde{\mathbf{w}} \right\} + \text{const.}$$

Design matrix and target vector:

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top \\ \tilde{\mathbf{x}}_2^\top \\ \vdots \\ \tilde{\mathbf{x}}_N^\top \end{pmatrix} \in \mathbb{R}^{N \times (D+1)}, \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \in \mathbb{R}^N$$

Compact expression:

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^\top \mathbf{y} \right)^\top \tilde{\mathbf{w}} \right\} + \text{const}$$

Example: $RSS(\tilde{\mathbf{w}})$ in Compact Form

| sqft (1000's) | bedrooms | bathrooms | sale price (100k) |
|---------------|----------|-----------|-------------------|
| 1 | 2 | 1 | 2 |
| 2 | 2 | 2 | 3.5 |
| 1.5 | 3 | 2 | 3 |
| 2.5 | 4 | 2.5 | 4.5 |

Design matrix and target vector:

$$\tilde{\mathbf{X}} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top \\ \tilde{\mathbf{x}}_2^\top \\ \vdots \\ \tilde{\mathbf{x}}_N^\top \end{pmatrix} = \begin{bmatrix} 1 & 1 & 2 & 1 \\ 1 & 2 & 2 & 2 \\ 1 & 1.5 & 3 & 2 \\ 1 & 2.5 & 4 & 2.5 \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} 2 \\ 3.5 \\ 3 \\ 4.5 \end{bmatrix}$$

. Compact expression:

$$RSS(\tilde{\mathbf{w}}) = \|\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}\|_2^2 = \left\{ \tilde{\mathbf{w}}^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^\top \mathbf{y} \right)^\top \tilde{\mathbf{w}} \right\} + \text{const}$$

Solution in Matrix Form

Compact expression

$$RSS(\tilde{\mathbf{w}}) = ||\tilde{\mathbf{X}}\tilde{\mathbf{w}} - \mathbf{y}||_2^2 = \left\{ \tilde{\mathbf{w}}^\top \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2 \left(\tilde{\mathbf{X}}^\top \mathbf{y} \right)^\top \tilde{\mathbf{w}} \right\} + \text{const}$$

Gradients of Linear and Quadratic Functions

- $\nabla_{\mathbf{x}}(\mathbf{b}^\top \mathbf{x}) = \mathbf{b}$
- $\nabla_{\mathbf{x}}(\mathbf{x}^\top \mathbf{A} \mathbf{x}) = 2\mathbf{A} \mathbf{x}$ (symmetric \mathbf{A})

Normal equation

$$\nabla_{\tilde{\mathbf{w}}} RSS(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}^\top \mathbf{y} = 0$$

This leads to the **least-mean-squares** (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{y}$$

Example: $RSS(\tilde{\mathbf{w}})$ in Compact Form

| sqft (1000's) | bedrooms | bathrooms | sale price (100k) |
|---------------|----------|-----------|-------------------|
| 1 | 2 | 1 | 2 |
| 2 | 2 | 2 | 3.5 |
| 1.5 | 3 | 2 | 3 |
| 2.5 | 4 | 2.5 | 4.5 |

Write the **least-mean-squares** (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^T \mathbf{y}$$

Can use solvers in Matlab, Python etc., to compute this for any given $\tilde{\mathbf{X}}$ and \mathbf{y} .

Exercise: $RSS(\tilde{\mathbf{w}})$ in Compact Form

Using the general **least-mean-squares** (LMS) solution

$$\tilde{\mathbf{w}}^{LMS} = \left(\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right)^{-1} \tilde{\mathbf{X}}^\top \mathbf{y}$$

recover the uni-variate solution that we had computed earlier:

$$w_1 = \frac{\sum (x_n - \bar{x})(y_n - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{and} \quad w_0 = \bar{y} - w_1 \bar{x}$$

where $\bar{x} = \frac{1}{N} \sum_n x_n$ and $\bar{y} = \frac{1}{N} \sum_n y_n$.

Exercise: $RSS(\tilde{\mathbf{w}})$ in Compact Form

For the 1-D case, the **least-mean-squares** solution is

$$\begin{aligned}\tilde{\mathbf{w}}^{LMS} &= (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^\top \mathbf{y} \\ &= \left(\begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ 1 & \dots \\ 1 & x_N \end{bmatrix} \right)^{-1} \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_N \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{bmatrix} \\ &= \left(\begin{bmatrix} N & N\bar{x} \\ N\bar{x} & \sum_n x_n^2 \end{bmatrix} \right)^{-1} \begin{bmatrix} \sum_n y_n \\ \sum_n x_n y_n \end{bmatrix} \\ \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} &= \frac{1}{\sum (x_i - \bar{x})^2} \begin{bmatrix} \bar{y} \sum (x_i - \bar{x})^2 - \bar{x} \sum (x_n - \bar{x})(y_n - \bar{y}) \\ \sum (x_n - \bar{x})(y_n - \bar{y}) \end{bmatrix}\end{aligned}$$

where $\bar{x} = \frac{1}{N} \sum_n x_n$ and $\bar{y} = \frac{1}{N} \sum_n y_n$.

Recap of MLE/MAP

Linear Algebra Review

Linear Regression

Formulation

Univariate Solution

Multivariate Solution

Probabilistic Interpretation

Why Minimize the RSS?

Probabilistic interpretation

- **Noisy observation model** for generating the dataset:

$$Y = w_0 + w_1 X + \eta$$

where $\eta \sim N(0, \sigma^2)$ is a Gaussian random variable

- Conditional likelihood of one training sample:

$$p(y_n|x_n) = N(w_0 + w_1 x_n, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{[y_n - (w_0 + w_1 x_n)]^2}{2\sigma^2}}$$

Probabilistic Interpretation (cont'd)

Log-likelihood of the training data \mathcal{D} (assuming i.i.d):

$$\begin{aligned}\log P(\mathcal{D}) &= \log \prod_{n=1}^N p(y_n|x_n) = \sum_n \log p(y_n|x_n) \\&= \sum_n \left\{ -\frac{[y_n - (w_0 + w_1 x_n)]^2}{2\sigma^2} - \log \sqrt{2\pi}\sigma \right\} \\&= -\frac{1}{2\sigma^2} \sum_n [y_n - (w_0 + w_1 x_n)]^2 - \frac{N}{2} \log \sigma^2 - N \log \sqrt{2\pi} \\&= -\frac{1}{2} \left\{ \frac{1}{\sigma^2} \sum_n [y_n - (w_0 + w_1 x_n)]^2 + N \log \sigma^2 \right\} + \text{const}\end{aligned}$$

What is the relationship between minimizing RSS and maximizing the log-likelihood?

Maximum Likelihood Estimation

$$\log P(\mathcal{D}) = -\frac{1}{2} \left\{ \frac{1}{\sigma^2} \sum_n [y_n - (w_0 + w_1 x_n)]^2 + N \log \sigma^2 \right\} + \text{const}$$

Estimating σ , w_0 and w_1 can be done in two steps

- Maximize over w_0 and w_1 :

$$\max \log P(\mathcal{D}) \Leftrightarrow \min \sum_n [y_n - (w_0 + w_1 x_n)]^2 \leftarrow \text{This is RSS}(\tilde{\mathbf{w}})!$$

- Maximize over $s = \sigma^2$:

$$\begin{aligned} \frac{\partial \log P(\mathcal{D})}{\partial s} &= -\frac{1}{2} \left\{ -\frac{1}{s^2} \sum_n [y_n - (w_0 + w_1 x_n)]^2 + N \frac{1}{s} \right\} = 0 \\ \rightarrow \sigma^{*2} = s^* &= \frac{1}{N} \sum_n [y_n - (w_0 + w_1 x_n)]^2 \end{aligned}$$

Why Is This Interpretation Useful?

- It gives a solid footing to our intuition: minimizing $\text{RSS}(\tilde{\mathbf{w}})$ is a sensible thing based on reasonable modeling assumptions.
- Estimating σ^* tells us how much noise there is in our predictions. For example, it allows us to place confidence intervals around our predictions.

You Should Know

- Linear regression is the linear combination of features
 $f : \mathbf{x} \rightarrow y$, with $f(\mathbf{x}) = w_0 + \sum_d w_d x_d = w_0 + \mathbf{w}^\top \mathbf{x}$
- If we minimize residual sum of squares as our learning objective, we get a closed-form solution of parameters
- Probabilistic interpretation: maximum likelihood if assuming residual is Gaussian distributed