

Read:

Murphy - 4.2, 7.1-7.3, 7.8, 11.1-11.2

Recall the summary of Maximum Likelihood Estimation (MLE).

Observe: x_1, x_2, \dots, x_n drawn iid from $f(x; \theta)$ for some "true" $\theta = \theta^*$

Likelihood function: $L_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$

Log-Likelihood function:

$$l_n(\theta) = \log(L_n(\theta)) = \sum_{i=1}^n \log(f(x_i; \theta))$$

MLE: $\hat{\theta}_{MLE} = \arg \max_{\theta} L_n(\theta)$

- As the number of observations $n \rightarrow \infty$, we have $\hat{\theta}_{MLE} \rightarrow \theta^*$

- MLE is a recipe that begins with a model for data $f(x; \theta)$

when Recipe is applied to continuous case:

Learning Gaussian Parameters (i.e., our model is assumed to be a Gaussian)

MLE: $\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

- MLE for the variance of a Gaussian is biased!

$$E[\hat{\sigma}_{MLE}^2] \neq \sigma^2 \quad (\text{see also definition of biased in Lecture \#1})$$

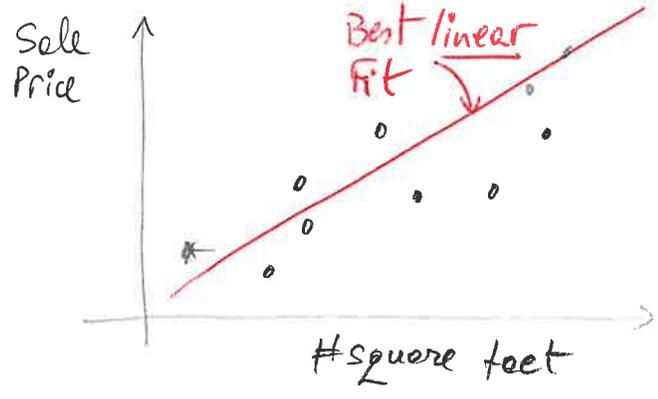
- Unbiased variance estimator:

$$\hat{\sigma}_{\text{unbiased}}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

Regression - 1-dimensional (simple linear regression)

Given past sales data on Zillow.com, predict:

y = house sale price from:
 $x = \{ \# \text{sq. ft.} \}$



Training Data: $\{(x_i, y_i)\}_{i=1}^n$ $x_i, y_i \in \mathbb{R}$

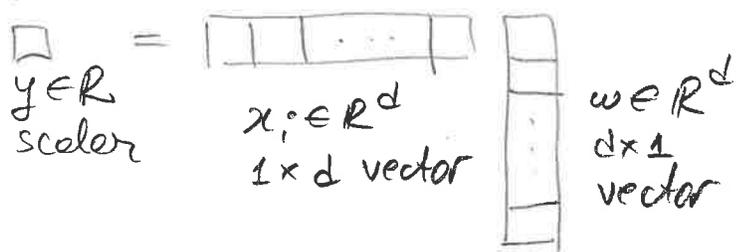
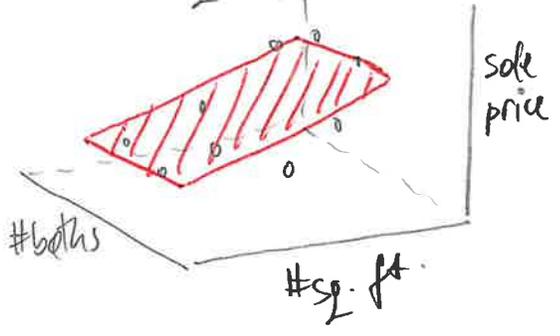
Hypothesis / Model: linear

$$y_i = x_i \cdot w + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$w \in \mathbb{R}$ is the slope of line weights or regression coefficients. **Gaussian Noise!**

d-dimensional (multiple linear regression): $x_i \in \mathbb{R}^d$

y = house sale price
 $x = \{ \# \text{sq. ft.}, \# \text{baths, etc.} \}$



Hypothesis / Model: Linear

$$y_i = x_i^T \cdot w + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$$P(y|x, w, \sigma) = N(x^T w, \sigma^2) \Rightarrow$$

$$P(y|x, w, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(y - x^T w)^2}{2\sigma^2}}$$

(Murphy pp. 110)

Now, maximize Log-Likelihood:

Training Data

$$\{(x_i, y_i)\}_{i=1}^n, x_i \in \mathbb{R}^d, y_i \in \mathbb{R} \quad p(y|x, \omega, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y-x^T\omega)}{2\sigma^2}}$$

Likelihood:

$$P(\mathcal{D}|\omega, \sigma) = \prod_{i=1}^n p(y_i|x_i, \omega, \sigma) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y_i-x_i^T\omega)}{2\sigma^2}}$$

Maximize - with respect to ω :

$$\begin{aligned} \log P(\mathcal{D}|\omega, \sigma) &= \log \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y_i-x_i^T\omega)}{2\sigma^2}} \right) \\ &= \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y_i-x_i^T\omega)}{2\sigma^2}} \right] \end{aligned}$$

use:
 $\log A \cdot B = \log A + \log B$

$$= n \cdot \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \sum_{i=1}^n \frac{(y_i - x_i^T \omega)^2}{2\sigma^2}$$

does not depend on ω

$$\hat{\omega}_{MLE} = \arg \max_{\omega} - \sum_{i=1}^n \frac{(y_i - x_i^T \omega)^2}{2\sigma^2}$$

constant when looking wrt ω

$$\hat{\omega}_{MLE} = \arg \min_{\omega} \sum_{i=1}^n (y_i - x_i^T \omega)^2$$

Minimize squared error!

The regression problem in matrix notation:

$$\hat{\omega}_{MLE} = \arg \min_{\omega} \sum_{i=1}^n (y_i - x_i^T \omega)^2$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Labels

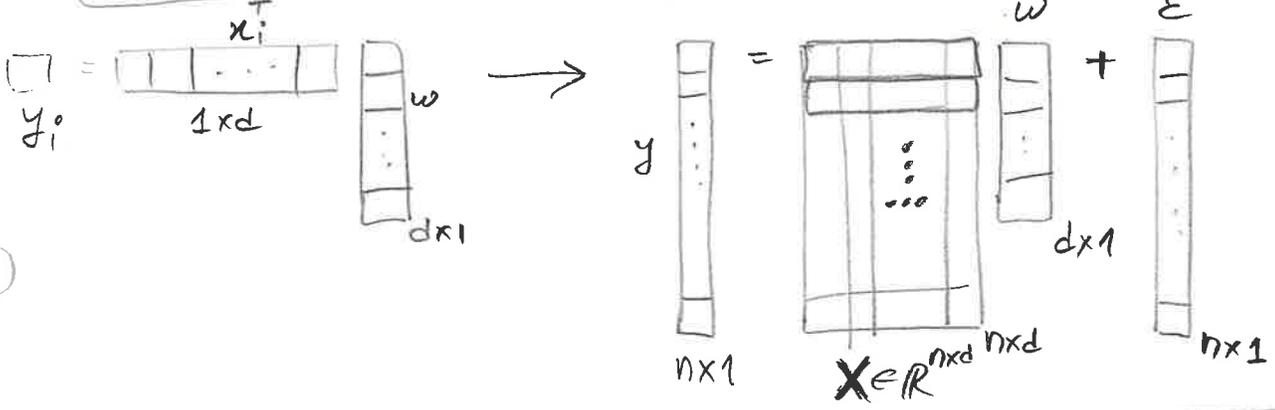
$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix}$$

Features

n is # of examples/data points
 d is # of features.

$$x_1^T = [x_{11}, x_{12}, \dots, x_{1d}]_{1 \times d}$$

matrix format $y_i = x_i^T \cdot \omega + \epsilon_i$ \rightarrow $y = X \cdot \omega + \epsilon$



l_2 norm: $\|z\|_2 = \sqrt{\sum_{i=1}^n z_i^2} \equiv \sqrt{z^T \cdot z}$ (Murphy pp. 229)

$$z_i = y_i - x_i^T \cdot \omega$$

Least Squares: (works with Matrices)

$$\hat{\omega}_{LS} = \arg \min_{\omega} (\|y - X \cdot \omega\|_2)^2$$

$$\hat{\omega}_{LS} = \arg \min_{\omega} (y - X \omega)^T \cdot (y - X \omega)$$

\triangleq Residual Sum of Squares (RSS)

[Murphy 11.2, pp. 367]

$$= \arg \min_w (y - Xw)^T (y - Xw)$$

$$= \arg \min_w y^T y - y^T Xw - (Xw)^T y + (Xw)^T (Xw)$$

Use:
 $(ABC)^T = C^T B^T A^T$
 $S \in \mathbb{R}, s = s^T$

$$= \arg \min_w \cancel{y^T y} - \underbrace{y^T X w}_{1 \times n \cdot n \times d \cdot d \times 1} - \underbrace{(w^T X^T y)^T}_{1 \times d \cdot d \times n \cdot n \times 1} + w^T X^T X w$$

Useful gradients:

- ① $\nabla_w (X^T w) = X$
- ② $\nabla_w (X^T A w) = A^T X$
- ③ $\nabla_w (w^T A w) = (A + A^T) w$
 if $A = A^T$ symmetric then:
 $\nabla_w (w^T A w) = 2Aw$

$$= \arg \min_w -y^T Xw - y^T Xw + w^T X^T X w$$

$$= \arg \min_w -2y^T Xw + w^T X^T X w$$

$$\nabla_w [-2y^T Xw + w^T X^T X w] = 0$$

$$-2X^T y + 2X^T X w = 0$$

multiply with $(X^T X)^{-1}$ | $X^T X w = X^T y$

NOTE: ① this assumes $X^T X$ is invertible!
 ② if it is singular (determinant = 0) => many optimal solutions exist!

$$\hat{w}_{LS} = (X^T X)^{-1} X^T y$$

← unique sol.

$$X^{\dagger} = (X^T X)^{-1} X^T$$

called the pseudo-inverse in Abu-Mustafa
 ...and would define X^{\dagger} in other ways if $X^T X$ is singular!

$$\hat{w}_{LS} = \left(\sum_{i=1}^n x_i x_i^T \right)^{-1} \sum_{i=1}^n x_i y_i$$

$\begin{matrix} dx_1 & dx_2 \\ dx_1 & dx_2 \end{matrix}$

Hence, we note that:

Linear Regression \equiv **OLS** (Ordinary) Least Squares

$$\hat{w}_{LS} = \hat{w}_{MLE} = (X^T X)^{-1} X^T y$$

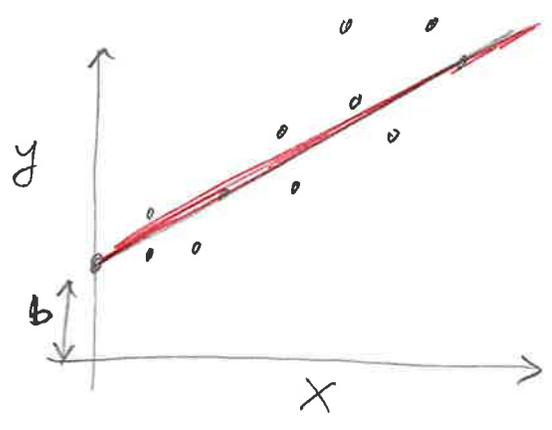
Called "Normal Equation" in Geron's Book (pp. 131)

What about an offset?

$$\hat{w}_{LS} = \arg \min_w \|y - Xw\|_2^2$$

$$= (X^T X)^{-1} X^T y$$

bias,
offset,
intercept.



$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \sum_{i=1}^n (y_i - (x_i^T w + b))^2$$

$$\hat{w}_{LS}, \hat{b}_{LS} = \arg \min_{w, b} \left\| y - (Xw + \mathbf{1}b) \right\|_2^2 \leftarrow \text{in matrix format.}$$

$n \times 1$ $b \in \mathbb{R}$
 scalar

$$\begin{cases} X^T X \hat{w}_{LS} + \hat{b}_{LS} X^T \mathbf{1} = X^T y \\ \mathbf{1}^T X \hat{w}_{LS} + \hat{b}_{LS} \mathbf{1}^T \mathbf{1} = \mathbf{1}^T y \end{cases}$$

$$\begin{bmatrix} X^T X & X^T \mathbf{1} \\ \mathbf{1}^T X & \mathbf{1}^T \mathbf{1} \end{bmatrix} \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} X^T y \\ \mathbf{1}^T y \end{bmatrix}$$

$$\tilde{x}_i = \begin{bmatrix} x_i \\ \mathbf{1} \end{bmatrix}$$

scalar $\in \mathbb{R}$

$$\tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$y_i = \tilde{x}_i^T \tilde{w}$$

Murphy introduces notation $b \triangleq w_0$, which is essentially absorbed in w .

$$\tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \tilde{x}_2^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} = [X \ \mathbf{1}]$$

$$y = \tilde{X}^T \tilde{w}$$

$$\hat{\tilde{w}}_{MLE} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$$

Appendix A

Derivation of the two-equations expressed as a compact block-matrix system \equiv is mathematically equivalent to augmenting X with a column of 1's and solving the usual normal equations.

- Note: augmenting X can be done with ones at the end or at the beginning

- Start with: $\hat{w}_{LS}, \hat{b}_{LS} = \underset{w, b}{\operatorname{argmin}} \left(\|y - (Xw + 1b)\|_2 \right)^2$; matrix format.
"L(theta)" - "loss function", $\theta = (w, b)$

- Derivative w.r.t. w :
gradient

$$\nabla_w L(w, b) = -2 X^T (y - Xw - \underbrace{1b}_{n \times 1 \text{ scalar}})$$

Set to zero:

$$X^T y = X^T X w + b X^T 1 \quad (1)$$

$$1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}_{n \times 1}$$

- Gradient w.r.t. b :

$$\nabla_b L(w, b) = -2 \cdot 1^T (y - Xw - 1b)$$

Set to zero:

$$1^T y = 1^T X w + b 1^T 1 \quad (2)$$

$$(1), (2) \Rightarrow \begin{bmatrix} X^T y \\ 1^T y \end{bmatrix} = \begin{bmatrix} X^T X & X^T 1 \\ 1^T X & 1^T 1 \end{bmatrix} \cdot \begin{bmatrix} w \\ b \end{bmatrix} \quad (3)$$

If this is invertible, then, we can derive:

$$\begin{bmatrix} \hat{w} \\ \hat{b} \end{bmatrix} = \begin{bmatrix} X^T X & X^T 1 \\ 1^T X & 1^T 1 \end{bmatrix}^{-1} \cdot \begin{bmatrix} X^T y \\ 1^T y \end{bmatrix}$$

which is equivalent to what we would get if we appended x_i with ones, append w with b , write matrix format and solve usual normal regression as done before:

$$\tilde{x}_i = \begin{bmatrix} x_i \\ 1 \end{bmatrix}, \tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix} \quad \tilde{y}_i = \tilde{x}_i^T \tilde{w} \xrightarrow{\text{matrix format}} y = \tilde{X}^T \tilde{w}, \text{ with } \tilde{X} = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = [X \ 1]$$

$$\text{Linear Regression} \Rightarrow \hat{w}_{MLE} = (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y$$

So, idea is to preprocess features to have mean zero.

if $X^T \cdot 1 = 0$ (i.e., if each feature is mean-zero) then:

$$\hat{\omega}_{LS} = (X^T X)^{-1} X \cdot y$$

$$\hat{b}_{LS} = \frac{1}{n} \sum_{i=1}^n y_i$$

Now, Making Predictions:

"A new house is about to be listed. What should it sell for?"

$$\hat{y}_{new} = x_{new}^T \hat{\omega}_{LS} + \hat{b}_{LS}$$

→ if train features were normalized have to normalize by subtracting the training mean.

Process:

- Decide on a model for the likelihood function $f(x; \theta)$
- Find the function that fits best the data
 - > Choose a loss function - least squares.
 - > Pick the function that minimizes loss on data.
- Use function to make prediction on new examples.

Recap so far

ML algorithms we looked at so far:

Maximum Likelihood Estimation (MLE)

x
 $p(x)$ { \rightarrow Fit a Bernoulli distribution (coin flips)
 \rightarrow Fit a Gaussian distribution (μ, σ)

x, y
 $p(y|x)$ \rightarrow Fit a Linear Predictor $x \rightarrow y$ with Gaussian Noise