

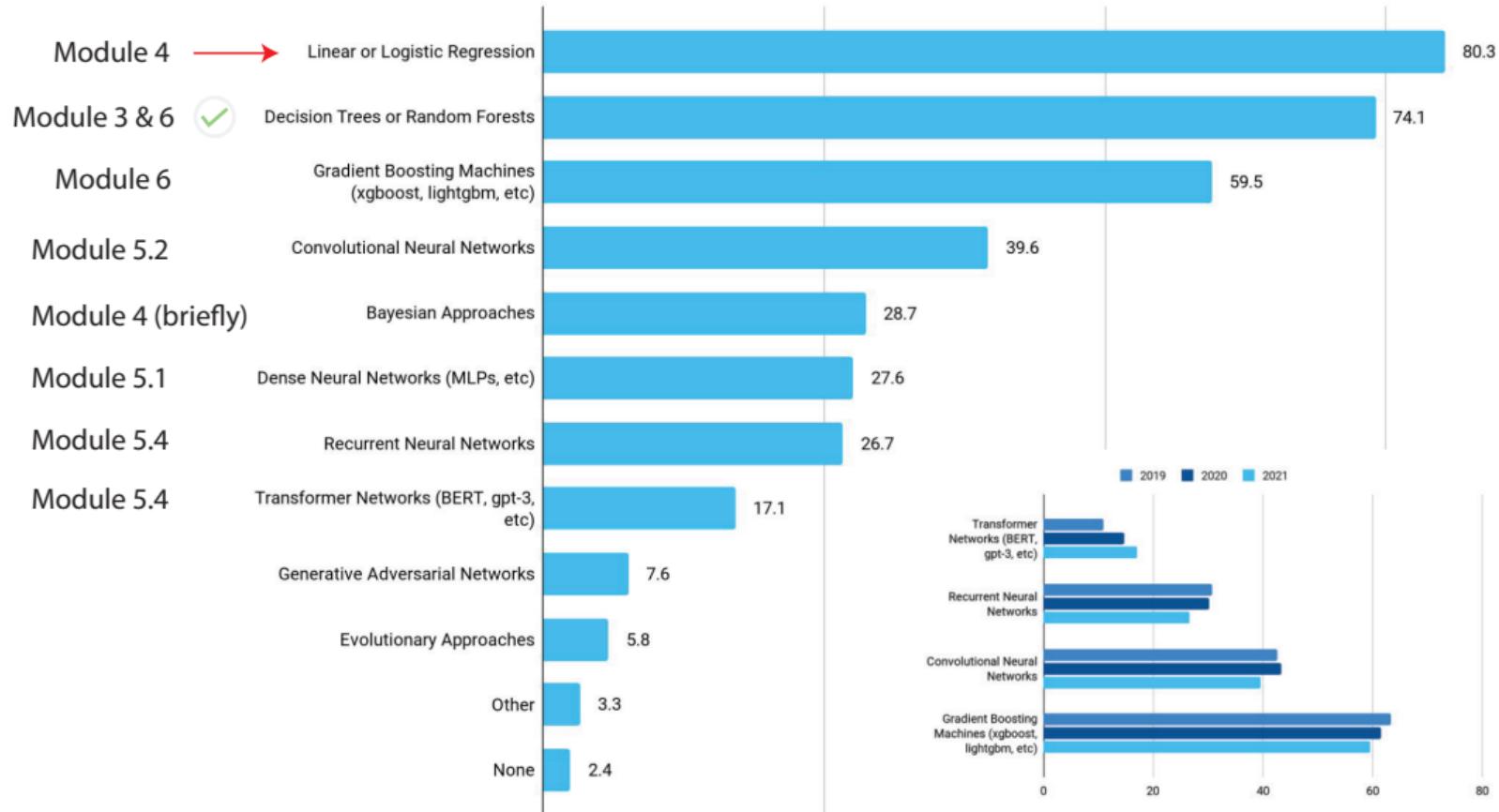
Lecture 7 & 8 - Module 4.1. Linear regression

COMP 551 Applied machine learning

Yue Li
Assistant Professor
School of Computer Science
McGill University

January 28 & 30, 2025

Most commonly used ML algorithms in Kaggle 2021 survey



Outline

Objectives

Simple linear regression

Multiple linear regression

Probabilistic interpretation

Non-linear basis function

Summary

Outline

Objectives

Simple linear regression

Multiple linear regression

Probabilistic interpretation

Non-linear basis function

Summary

Learning objectives

Understanding the following concepts

- Simple linear model
- Finding the best linear fit by minimizing SSE
- Matrix algebra in solving multiple regression
- Probabilistic interpretation
- Feature transformation by non-linear basis functions

Outline

Objectives

Simple linear regression

Multiple linear regression

Probabilistic interpretation

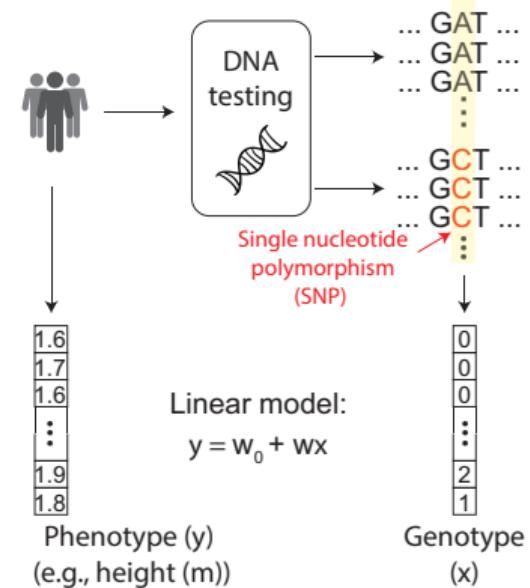
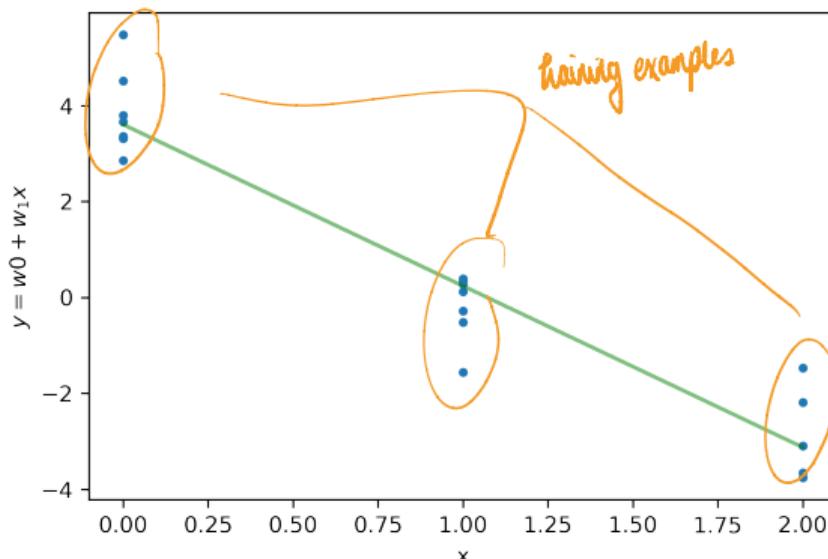
Non-linear basis function

Summary

Linear regression using a one-dimensional input

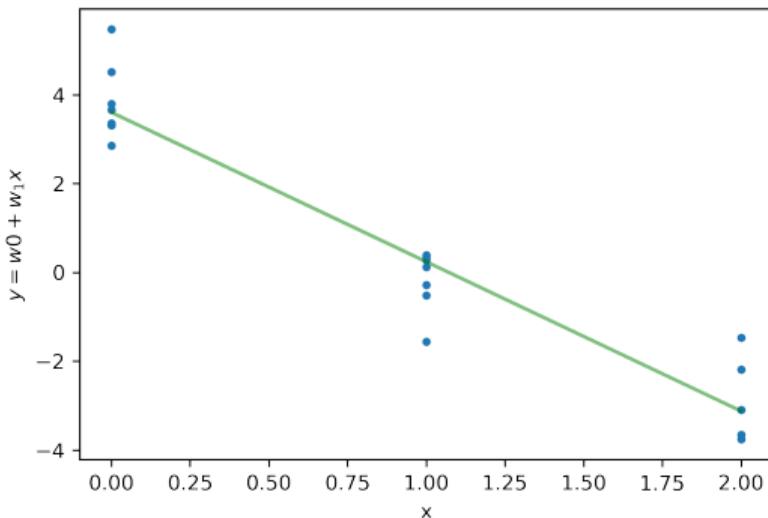
We want to predict real-valued quantity or often known as **response or target variable** $y \in \mathbb{R}$ by finding a mapping function that maps from a **one-dimensional input** x to the real-valued y . We can fit a linear function on training examples $\{x^{(n)}, y^{(n)}\}_{n=1}^N$ (e.g., predicting phenotype from one genetic mutation):

$$f(x^{(n)}; w_0, w_1) = w_0 + w_1 x^{(n)} \quad (1)$$



Linear regression using a one-dimensional input

$$f(x^{(n)}; w_0, w_1) = w_0 + w_1 x^{(n)}$$



- $w_0 = 3.4$ is the **intercept** or **bias**, which is not be confused with the “model bias”
- $w_1 = -3.25$ is the **slope** of the linear function or the **regression coefficient**.

Simple linear regression using one input feature

In statistics, we often write down the regression formula as:

$$\epsilon^{(n)} = \hat{y}^{(n)} - \bar{y}^{(n)}$$

$$y^{(n)} = w_0 + w_1 x^{(n)} + \epsilon^{(n)}$$

predicted value $\rightarrow \hat{y}^{(n)}$ *prediction error*

where $\epsilon^{(n)}$ is the prediction error for example n .

Using matrix notation, we can write the regression formula as:

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix}_{N \times 1} = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x^{(2)} \\ \vdots & \vdots \\ 1 & x^{(N)} \end{bmatrix}_{N \times 2} \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}_{2 \times 1} + \begin{bmatrix} \epsilon^{(1)} \\ \epsilon^{(2)} \\ \vdots \\ \epsilon^{(N)} \end{bmatrix}_{N \times 1} \quad (2)$$

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \boldsymbol{\epsilon} \quad (3)$$

where \mathbf{y} and $\boldsymbol{\epsilon}$ are both $N \times 1$ vectors, \mathbf{X} is a $N \times 2$ matrix, and \mathbf{w} is a 2×1 vector.

Residual error as a measure of prediction loss

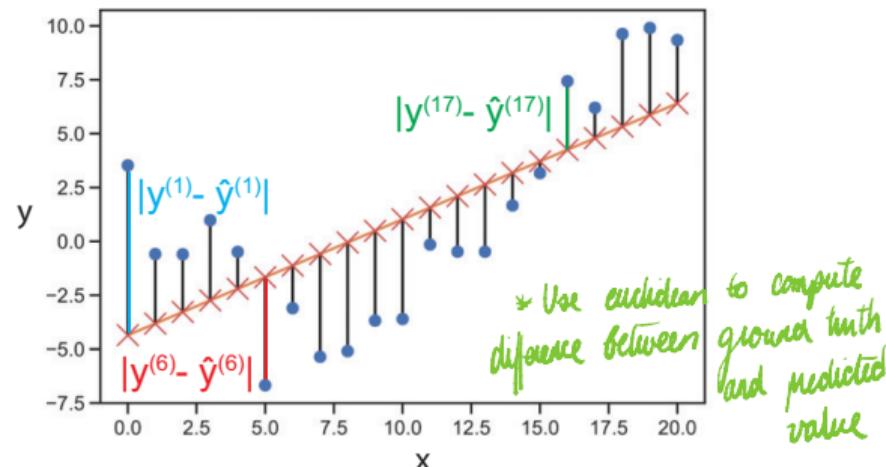
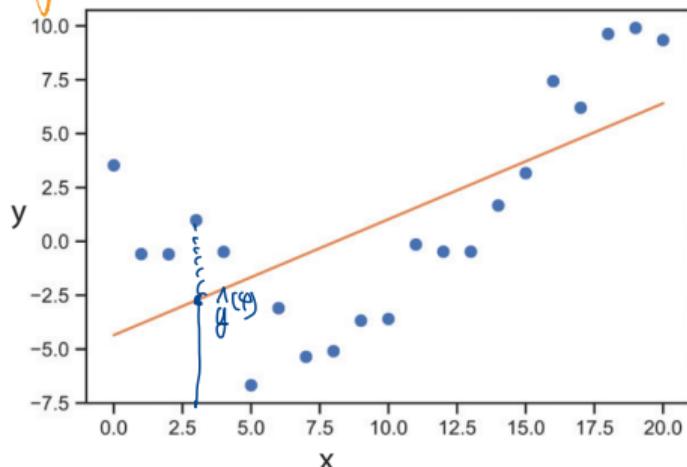
Every straight line we fit incurs a prediction error on the training data point unless the fitted line goes through that data point. The **residual error** is Euclidean distance between the observed response $y^{(n)}$ value and the predicted response $\hat{y}^{(n)} = \mathbf{x}^{(n)}\mathbf{w}$:

loss function $I_n = \|\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)}\|_2 = \sqrt{(\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)})^2} = |\mathbf{y}^{(n)} - \hat{\mathbf{y}}^{(n)}|$ (4)

Fitting a linear function on a 1D input and output data

$\sqrt{(\varepsilon^{(n)})^2} \geq 0$ $|\varepsilon^{(n)}|$

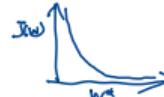
Error from the true value



Notation: $\|\mathbf{a}\|_2 = \sqrt{\sum_i a_i^2}$ is L₂-norm and $\|\mathbf{a}\|_2^2 = \sum_i a_i^2$ is the squared of the L2-norm.

Fitting a linear regression function by minimizing the sum of squared error

Sum of squared error (SSE) as a function of the linear coefficients \mathbf{w} is defined as:

$$J(\mathbf{w}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \sum_{n=1}^N (y^{(n)} - \hat{y}^{(n)})^2 = \sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)})^2 = (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{e}^T \mathbf{e} = \sum_n (\epsilon^{(n)})^2$$


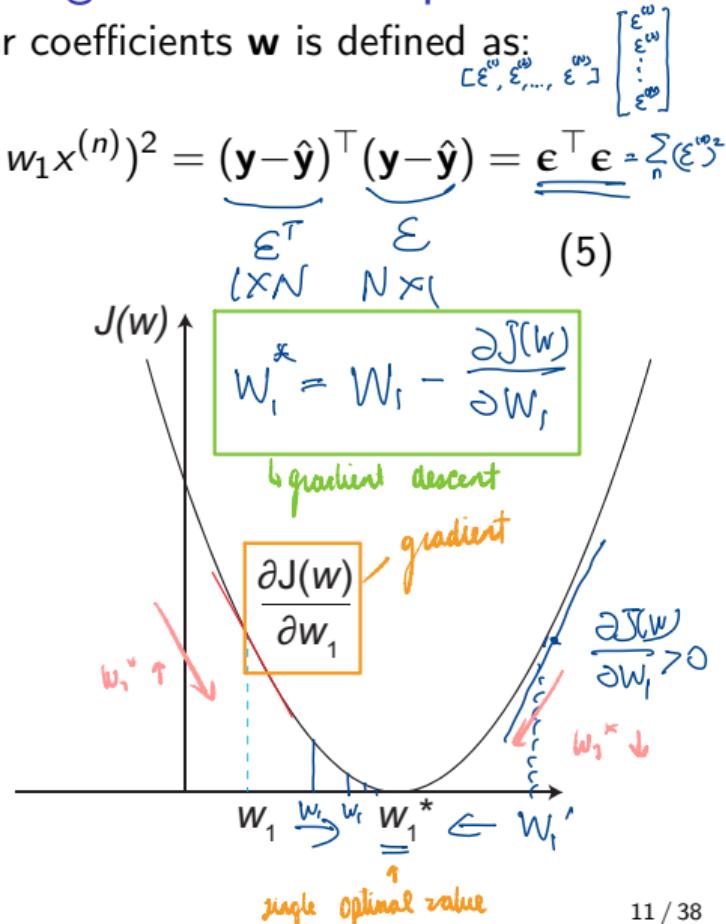
Goal: find the best \mathbf{w} that minimizes the SSE:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w}) \quad (6)$$

*Key for optimization

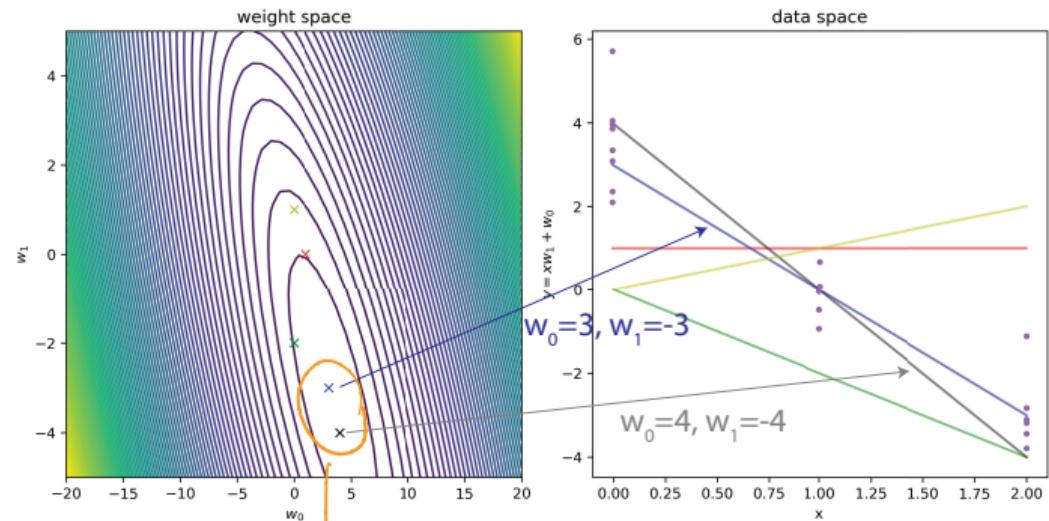
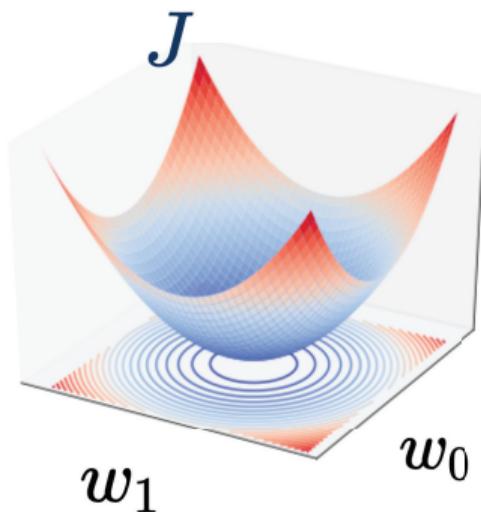
Given the error function is differentiable everywhere, the slope at the current value of w_1 projecting onto the error parabola is shown on the right.

The SSE error function is **convex** (more details in Module 4.4). The optimal w_1^* is found where **the slope or the partial derivative is zero** $\frac{\partial J(\mathbf{w})}{\partial w_1^*} = 0$.



Contour plot visualizes the error surface as a function of w_1 and w_0

$$J(\mathbf{w}) = \sum_{n=1}^N (y^{(n)} - \hat{y}^{(n)})^2 = \sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)})^2$$



center is optimal value

Derivation of ordinary least squared (OLS) solution for w_0

$$\frac{\partial \alpha^2}{\partial a} = 2a$$

$$\frac{\partial -w_0}{\partial w_0} = -1$$

Solving for w_0 :

$$\sum_{n=1}^N y^{(n)} - \sum_{n=1}^N w_0 - \sum_{n=1}^N w_1 x^{(n)} = 0$$

$$\sum_{n=1}^N w_0 = \sum_{n=1}^N y^{(n)} - w_1 \sum_{n=1}^N x^{(n)}$$

$$\underline{Nw_0} = \sum_{n=1}^N y^{(n)} - w_1 \sum_{n=1}^N x^{(n)}$$

w_0 never changes

$$w_0 = \frac{1}{N} \sum_{n=1}^N y^{(n)} - w_1 \frac{1}{N} \sum_{n=1}^N x^{(n)}$$

$$w_0 = \bar{y} - w_1 \bar{x}$$

any input value
any response value

Therefore, $\hat{w}_0 = \bar{y} - w_1 \bar{x}$

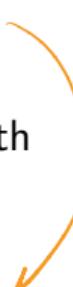
$$\frac{\partial J(w)}{\partial w_0} = \frac{\partial}{\partial w_0} \sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)})^2$$

$$= \sum_{n=1}^N \frac{\partial (y^{(n)} - w_0 - w_1 x^{(n)})^2}{\partial (y^{(n)} - w_0 - w_1 x^{(n)})} \frac{\partial (y^{(n)} - w_0 - w_1 x^{(n)})}{\partial w_0}$$

$$= \sum_{n=1}^N 2(y^{(n)} - w_0 - w_1 x^{(n)})(-1) = 0$$

Set $\frac{\partial J(w_0)}{\partial w_0}$ to zero and multiplying $\frac{1}{2}$ and -1 on both sides gives:

$$\sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)}) = 0$$



Derivation of OLS solution for w_1

$$J = \sum_{n=1}^N (y^{(n)} - \bar{y})^2 = \sum_n (\varepsilon^{(n)})^2 = (\varepsilon^{(1)})^2 + (\varepsilon^{(2)})^2 + \dots + (\varepsilon^{(N)})^2$$

$$\frac{\partial J(\mathbf{w})}{\partial w_1} = \frac{\partial}{\partial w_1} \sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)})^2$$

$$= \sum_{n=1}^N 2(y^{(n)} - w_0 - w_1 x^{(n)})(-x^{(n)}) = 0$$

Set $\frac{\partial J(\mathbf{w})}{\partial w_1}$ to zero and multiplying $\frac{1}{2}$ and -1 on both sides gives:

$$\sum_{n=1}^N (y^{(n)} - w_0 - w_1 x^{(n)})x^{(n)} = 0 \quad (7)$$

containing both response
and input

Plug the bias estimate \hat{w}_0 in (7) and solve for w_1 :

$$\sum_{n=1}^N (y^{(n)} - (\bar{y} - w_1 \bar{x}) - w_1 x^{(n)})x^{(n)} = 0$$

$$\sum_{n=1}^N [(y^{(n)} - \bar{y}) + (w_1 \bar{x} - w_1 x^{(n)})]x^{(n)} = 0$$

$$\sum_{n=1}^N (y^{(n)} - \bar{y})x^{(n)} - w_1 \sum_{n=1}^N (x^{(n)} - \bar{x})x^{(n)} = 0$$

$\hat{w}_1 = \frac{\sum_{n=1}^N (y^{(n)} - \bar{y})x^{(n)}}{\sum_{n=1}^N (x^{(n)} - \bar{x})x^{(n)}}$

→ Covariance ↑, \hat{w}_1 ↑

Covariance between x and y

Derivation of OLS solution for w_1 (cont'd)

Note that

Covariance of X and Y :

$$\text{Cov}(X, Y) = E[(Y - E[Y])(X - E[X])]$$

$$\sum_{n=1}^N (y^{(n)} - \bar{y})(x^{(n)} - \bar{x}) = \sum_{n=1}^N (y^{(n)} - \bar{y})(x^{(n)} - \bar{x}) \quad \text{Sample covariance between } X \text{ and } Y : \\ \bar{y} = \frac{1}{N} \sum_n y^{(n)} \quad N\bar{y} = \sum_n y^{(n)}$$

because:

$$\begin{aligned} \sum_{n=1}^N (y^{(n)} - \bar{y})(x^{(n)} - \bar{x}) &= \sum_{n=1}^N y^{(n)}x^{(n)} - \sum_{n=1}^N y^{(n)}\bar{x} - \sum_{n=1}^N \bar{y}x^{(n)} + \sum_{n=1}^N \bar{y}\bar{x} \\ &= \sum_{n=1}^N y^{(n)}x^{(n)} - N\bar{y}\bar{x} - \sum_{n=1}^N \bar{y}x^{(n)} + N\bar{y}\bar{x} = \sum_{n=1}^N (y^{(n)} - \bar{y})x^{(n)} \end{aligned}$$

Sample Variance :

$$S_{xx} = \frac{1}{N} \sum_n (x^{(n)} - \bar{x})^2$$

Similarly,

$$\sum_{n=1}^N (x^{(n)} - \bar{x})x^{(n)} = \sum_{n=1}^N (x^{(n)} - \bar{x})(x^{(n)} - \bar{x}) = \sum_{n=1}^N (x^{(n)} - \bar{x})^2$$

Update equations for simple linear regression

Therefore, the simple linear regression OLS solutions are:

$$\hat{w}_1 = \frac{\sum_{n=1}^N (y^{(n)} - \bar{y})x^{(n)}}{\sum_{n=1}^N (x^{(n)} - \bar{x})x^{(n)}} = \frac{\sum_{n=1}^N (y^{(n)} - \bar{y})(x^{(n)} - \bar{x})}{\sum_{n=1}^N (x^{(n)} - \bar{x})^2}; \quad \hat{w}_0 = \bar{y} - \hat{w}_1 \bar{x} \quad (8)$$

$$\hat{w}_1 = \frac{s_{xy}}{s_{xx}}$$

Compare this solution with Pearson correlation coefficient (PCC) yields useful insights:

$$\hat{r}_{xy} = \frac{\vec{y}^T \vec{x}}{\|\vec{x}\|_2 \|\vec{y}\|_2}$$

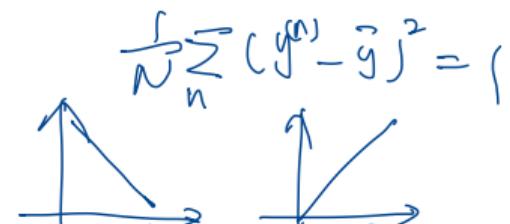
$$\hat{r}_{xy} = \frac{\sum_{n=1}^N (y^{(n)} - \bar{y})(x^{(n)} - \bar{x})}{\sqrt{\sum_{n=1}^N (x^{(n)} - \bar{x})^2} \sqrt{\sum_{n=1}^N (y^{(n)} - \bar{y})^2}} \in [-1, 1]$$

$$\hat{r}_{xy} = \frac{s_{xy}}{\sqrt{s_{xx}} \sqrt{s_{yy}}}$$

Therefore, we can see that

$$\begin{aligned} \hat{y} &= \hat{w}_1 x + \hat{w}_0 \\ &\Rightarrow \hat{w}_1 x + \bar{y} - \hat{w}_1 \bar{x} = \hat{w}_1 (x - \bar{x}) + \bar{y} \\ &\Rightarrow \hat{w}_1 x + \bar{y} \end{aligned}$$

$$\hat{w}_1 = \hat{r}_{xy} \frac{\sqrt{\sum_{n=1}^N (y^{(n)} - \bar{y})^2}}{\sqrt{\sum_{n=1}^N (x^{(n)} - \bar{x})^2}}$$



If the standard deviation for y and x are the same (e.g., through standardization described next), $\hat{w}_1 = \hat{r}_{xy}$. However, in general, the regression slope and PCC are *not* the same: while \hat{r}_{xy} gives a bounded measure independent of the scale of the two variables, \hat{w}_1 measures the change in the expected value of y corresponding to 1-unit increase/decrease of x .

Update equations for simple linear regression

A special case arises when y and x are *centered*: $\dot{y}^{(n)} = y^{(n)} - \bar{y}$ and $\dot{x}^{(n)} = x^{(n)} - \bar{x}$, such that $\bar{\dot{y}} = 0$ and $\bar{\dot{x}} = 0$ because

$$\bar{\dot{y}} = \frac{1}{N} \sum_n \dot{y}^{(n)} = \frac{1}{N} \sum_n (y^{(n)} - \bar{y}) = \frac{1}{N} \sum_n y^{(n)} - \frac{1}{N} \sum_n \bar{y} = \bar{y} - \bar{y} = 0$$

As a result,

$$\hat{w}_1 = \frac{\dot{y}^T \dot{x}}{\dot{x}^T \dot{x}}$$

$$\hat{w}_1 = \frac{\sum_{n=1}^N (\dot{y}^{(n)} - \bar{\dot{y}})(\dot{x}^{(n)} - \bar{\dot{x}})}{\sum_{n=1}^N (\dot{x}^{(n)} - \bar{\dot{x}})^2} = \frac{\sum_{n=1}^N \dot{y}^{(n)} \dot{x}^{(n)}}{\underbrace{\sum_{n=1}^N (\dot{x}^{(n)})^2}}; \quad \hat{w}_0 = \bar{y} - \hat{w}_1 \bar{\dot{x}} = 0$$

Without the intercept, the linear function becomes $\hat{y}^{(n)} = w_1 \dot{x}^{(n)}$



If we only center the input variable such that $\bar{x} = 0$ but not y , then $w_0 = \bar{y}$ and the linear function $\hat{y}^{(n)} = w_1 \dot{x}^{(n)} + \bar{y}$ measures the *change* from the average \bar{y} .

Also, PCC of centered x and y is the same as cosine similarity between x and y .

Standardization involves centering and scaling the input feature

Another special case arises if we further scale \dot{x} by its standard deviation

$$\sigma_{\dot{x}} = \sqrt{\frac{1}{N} \sum_n (\dot{x}^{(n)})^2}$$

$$\tilde{x}^{(n)} = \dot{x}^{(n)} / \sigma_{\dot{x}}$$

then we have

$$\sigma_{\tilde{x}}^2 = \frac{1}{N} \sum_n (\tilde{x}^{(n)})^2 = \frac{1}{N} \sum_n (\dot{x}^{(n)} / \sigma_{\dot{x}})^2 = \frac{1}{N} \sum_n (\dot{x}^{(n)})^2 / \sigma_{\dot{x}}^2 = \sigma_{\dot{x}}^2 / \sigma_{\dot{x}}^2 = 1$$

The regression coefficient becomes more simplified (while $w_0 = 0$ after centering x, y):

$$\hat{w}_1 = \frac{\sum_{n=1}^N \tilde{y}^{(n)} \tilde{x}^{(n)}}{\sum_{n=1}^N (\tilde{x}^{(n)})^2} = \frac{1}{N} \sum_{n=1}^N \tilde{y}^{(n)} \tilde{x}^{(n)} = \frac{1}{N} \tilde{\mathbf{x}}^T \tilde{\mathbf{y}}$$

What's PCC between \tilde{x} and \tilde{y}
 $N \times N \times 1$ after standardization? $\frac{1}{N} \tilde{\mathbf{x}}^T \tilde{\mathbf{y}}$

Together, the procedure of centering and scaling of the response and/or the input features is common practice known as **standardization** mainly to simplify computation and to make the model robust to different numerical scales of the input and response.

Outline

Objectives

Simple linear regression

Multiple linear regression

Probabilistic interpretation

Non-linear basis function

Summary

Multiple linear regression

- Goal: learn how to fit a *multiple regression* to predict the outcome variable y using *multiple* input features
- We can write the response as a **weighted linear sum** of the input features:

$$\hat{y}^{(n)} = f(\mathbf{x}^{(n)}; \mathbf{w}) = w_0 + w_1 x_1^{(n)} + \underbrace{w_2 x_2^{(n)}}_{\text{everything computed}} + \dots + w_D x_D^{(n)} = w_0 + \sum_{d=1}^D w_d x_d^{(n)}$$

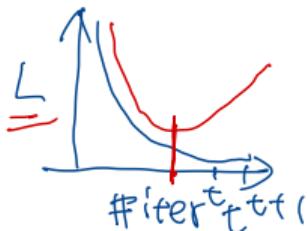
$$= \hat{w}_0 + \hat{w}_1 x_1^{(n)} + \hat{w}_2 x_2^{(n)} + \dots + \hat{w}_D x_D^{(n)} = \boxed{\hat{w}_0 + \sum_{d=2}^D \hat{w}_d x_d^{(n)}} + w_2 x_2$$

$= \hat{y}_{x_2} + \underbrace{w_2 x_2}$
fitting w_2

N individuals	... GAT AGG ...	1.6	0 1 0	...	0 0 0 0	?	\downarrow prediction without w_2	w_0
	... GAT AGG ...	1.7	0 0 1	...	2 0 2 0	?	w_0	
	... GAT AGG ...	1.6	1 0 1	...	0 1 0 1	?	w_0	
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
	... GCT ATG ...	1.9	0 0 0	...	0 0 0 0	?	w_0	
	... GCT ATG ...	1.8	0 1 1	...	1 0 1 0	?	w_0	
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	
	... GCT ATG ...	y phenotype	$N \times 1$	X genotype	$N \times D$	W coefficients	$D \times 1$	w_0 bias

Coordinate descent

$$\begin{aligned} \mathcal{E} &= y - \hat{y} \\ \underbrace{\mathcal{E}}_{N \times 1} &= y - \sum_{d=1}^D w_d x_d \end{aligned}$$



$$\begin{aligned} &= y - \underbrace{\sum_{d' \neq d} \hat{w}_{d'} x_{d'}}_{\Delta y_d} - w_d x_d \\ &= \Delta y_d - w_d x_d \end{aligned}$$

$$w_d^* \leftarrow \underset{w_d}{\operatorname{argmin}} \|y - \hat{y}\|_2^2 = \sum_{n=1}^N (\varepsilon^{(n)})^2$$

$$w_d^* = \frac{1}{N} \sum_{d=1}^D \Delta y_d$$

1. Initialize $W \in \mathbb{R}^{D \times 1}$
for $t = 1, \dots, T$:
2. for $d = 1, \dots, D$

$$\begin{aligned} \mathcal{E} &= \Delta y_d - w_d x_d \\ \text{Given } D-1 \quad \hat{w}_{d'}, \quad d' \neq d \end{aligned}$$

$$w_d^* = \frac{1}{N} \sum_{d=1}^D \Delta y_d$$

$$3. \quad \stackrel{(t)}{L} = \|y - \hat{y}\|_2^2$$

$$4. \quad \Delta L = \stackrel{(t)}{L} - \stackrel{(t-1)}{L}$$

threshold

if $\Delta L \leq T$, $T = 10^{-6}$
break

Multiple regression in matrix form

Suppose N training examples and D features. The data matrices are:

- Response: $\mathbf{y} \in \mathbb{R}^{N \times 1}$
- Input feature matrix: $[\mathbf{1}, \mathbf{X}] \in \mathbb{R}^{N \times (D+1)}$
- Regression coefficients: $[w_0; \mathbf{w}] \in \mathbb{R}^{(D+1) \times 1}$

We can rewrite the multiple regression function as:

$$\begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(N)} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_D^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{bmatrix} + \begin{bmatrix} \epsilon^{(1)} \\ \epsilon^{(2)} \\ \vdots \\ \epsilon^{(N)} \end{bmatrix}$$

$$\mathbf{y} = \underbrace{\mathbf{1} w_0}_{N \times 1} + \underbrace{\mathbf{X} \mathbf{w}}_{N \times 1} + \underbrace{\boldsymbol{\epsilon}}_{N \times 1}$$

Our goal is to find the ordinary least square (OLS) solution for the coefficients \mathbf{w} :

$$w_0^*, \mathbf{w}^* \leftarrow \arg \min_{w_0, \mathbf{w}} \|\mathbf{y} - \mathbf{X} \mathbf{w} - w_0\|_2^2 = (\mathbf{y} - \mathbf{X} \mathbf{w} - w_0)^T (\mathbf{y} - \mathbf{X} \mathbf{w} - w_0) = \sum_n (\epsilon^{(n)})^2$$

Notation: $\|\mathbf{a}\|_2 = \sqrt{\sum_i a_i^2}$ is L_2 -norm and $\|\mathbf{a}\|_2^2 = \sum_i a_i^2$ is the squared of the L2-norm.

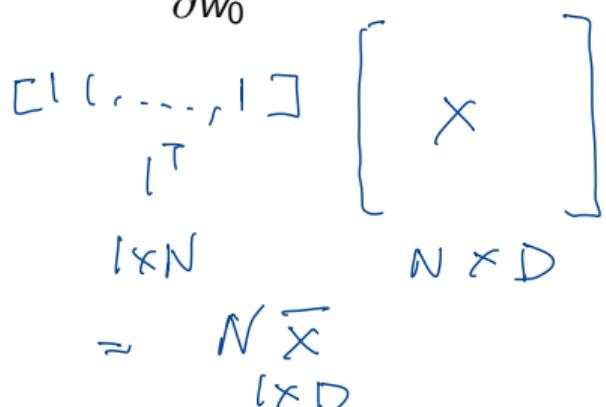
Multiple regression OLS derivation for the bias term w_0

Let the loss function be $J(\mathbf{w}, w_0) = (\mathbf{y} - \mathbf{X}\mathbf{w} - w_0)^T(\mathbf{y} - \mathbf{X}\mathbf{w} - w_0)$.

$$\begin{aligned}
 \frac{\partial J(\mathbf{w})}{\partial w_0} &= \frac{\partial}{\partial w_0} ((\mathbf{y} - \mathbf{X}\mathbf{w}) - \mathbf{1}w_0)^T((\mathbf{y} - \mathbf{X}\mathbf{w}) - \mathbf{1}w_0) \\
 &= \frac{\partial}{\partial w_0} \{(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) - (\mathbf{y} - \mathbf{X}\mathbf{w})^T\mathbf{1}w_0 - w_0\mathbf{1}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \mathbf{1}^T\mathbf{1}w_0^2\} \\
 &= \frac{\partial}{\partial w_0}(\mathbf{y} - \mathbf{X}\mathbf{w})^T(\mathbf{y} - \mathbf{X}\mathbf{w}) - \frac{\partial}{\partial w_0}2w_0\mathbf{1}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + \frac{\partial}{\partial w_0}Nw_0^2 \\
 &= 0 - 2\mathbf{1}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) + 2Nw_0 \\
 &= -2\mathbf{1}^T\mathbf{y} + 2\mathbf{1}^T\mathbf{X}\mathbf{w} + 2Nw_0 \\
 &= -2 \sum_n y^{(n)} + 2 \underbrace{\left(\sum_n \mathbf{x}^{(n)}\right)}_{\substack{1 \times D \\ \mathbf{X}}} \mathbf{w} + 2Nw_0 \stackrel{\text{set } 0}{=} 0
 \end{aligned}$$

Solving $\frac{\partial J(\mathbf{w})}{\partial w_0} = 0$ for w_0 :

$$\hat{w}_0 = \frac{1}{N} \sum_n y^{(n)} - \frac{1}{N} \left(\sum_n \mathbf{x}^{(n)} \right) \mathbf{w} \equiv \bar{y} - \bar{\mathbf{x}}\mathbf{w}$$



Plugging in the OLS estimate for the bias term w_0 into the loss

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{1}\bar{\mathbf{y}} \end{bmatrix} - \begin{bmatrix} \bar{\mathbf{y}} \\ \bar{\mathbf{y}} \\ \vdots \\ \bar{\mathbf{y}} \end{bmatrix}$$

$$\begin{aligned} J(\mathbf{w}, \hat{w}_0) &= \|\mathbf{y} - \mathbf{Xw} - \mathbf{1}\hat{w}_0\|_2^2 \\ &= \|\mathbf{y} - \mathbf{Xw} - \mathbf{1}(\bar{\mathbf{y}} - \bar{\mathbf{x}}\mathbf{w})\|_2^2 \\ &= \|(\mathbf{y} - \mathbf{1}\bar{\mathbf{y}}) - (\mathbf{Xw} - \mathbf{1}\bar{\mathbf{x}}\mathbf{w})\|_2^2 \\ &= \|(\mathbf{y} - \mathbf{1}\bar{\mathbf{y}}) - (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}})\mathbf{w}\|_2^2 \\ &\equiv \|\dot{\mathbf{y}} - \dot{\mathbf{Xw}}\|_2^2 \end{aligned}$$

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{1} \\ \bar{\mathbf{x}} \end{bmatrix} - \begin{bmatrix} \bar{\mathbf{x}} \\ \bar{\mathbf{x}} \\ \vdots \\ \bar{\mathbf{x}} \end{bmatrix}$$

where both the response and input features are mean-centered such that for each example n :

$$\dot{y}^{(n)} = y^{(n)} - \bar{y}$$

$$\dot{\mathbf{x}}^{(n)} = \mathbf{x}^{(n)} - \bar{\mathbf{x}}$$

Note: the second equation is element-wise subtraction: for each feature d ,

$$\dot{x}_d^{(n)} = x_d^{(n)} - \bar{x}_d$$

Solving all multiple regression coefficients simultaneously

Here we can assume that either the data are centered (i.e., $\mathbf{y} = \hat{\mathbf{y}}$, $\mathbf{X} = \hat{\mathbf{X}}$) or we have $\mathbf{x}_0 = \mathbf{1}$. In the latter case, the estimate for w_0 is the first element in $\hat{\mathbf{w}}$.

$$\begin{aligned}\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = \frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^\top \mathbf{y} - \underbrace{\mathbf{y}^\top \mathbf{X}\mathbf{w} - \mathbf{w}^\top \mathbf{X}^\top \mathbf{y}}_{2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y}} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}) \\ &= \frac{\partial}{\partial \mathbf{w}} (\mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}) = \underbrace{\frac{\partial}{\partial \mathbf{w}} \mathbf{y}^\top \mathbf{y}}_0 - 2 \underbrace{\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^\top \mathbf{X}^\top \mathbf{y}}_{\mathbf{b}^\top \mathbf{a}} + \underbrace{\frac{\partial}{\partial \mathbf{w}} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}}_{\mathbf{b}^\top \mathbf{A} \mathbf{b}} \\ &\stackrel{\dagger}{=} -2\mathbf{X}^\top \mathbf{y} + 2\underbrace{\mathbf{X}^\top \mathbf{X}\mathbf{w}}_{\mathbf{\Sigma}} = 2\mathbf{X}^\top (\underbrace{\mathbf{X}\mathbf{w} - \mathbf{y}}_{\mathbf{\Sigma}}) = 2\mathbf{X}^\top (\hat{\mathbf{y}} - \mathbf{y})\end{aligned}$$

To get this equality, we make use of two general properties in matrix differentiation¹:

$$\frac{\partial \mathbf{b}^\top \mathbf{a}}{\partial \mathbf{b}} = \mathbf{a}, \quad \frac{\partial \mathbf{b}^\top \mathbf{A}\mathbf{b}}{\partial \mathbf{b}} = 2\mathbf{A}\mathbf{b} \quad \frac{1}{N} \mathbf{X}_{\mathbf{d}}^\top \mathbf{X}_{\mathbf{d}}$$

Setting the derivative to zero and solve for \mathbf{w} gives the closed-form solution:

$$0 = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{w} \rightarrow \mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{y} \rightarrow \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

¹Matrix cookbook Section 2.4 Eq 69 & Eq 85

$D \times N$ $N \times D$

X, Y
 $N \times D \quad N \times 1$

$$\hat{W} = (X^T X)^{-1} X^T Y$$

$$= (N R)^{-1} N \hat{\beta}$$

$$= \frac{1}{N} R^{-1} N \hat{\beta}$$

$$= R^{-1} \hat{\beta}$$

$$\hat{y} = X^* \hat{W}$$

Summary Statistics (Sumstat)
Simple regression weights:

$$\hat{\beta} = \frac{1}{N} \tilde{X}^T \tilde{y}$$

PCC: $D \times 1 \quad D \times N \quad N \times 1$

$$R = \frac{1}{N} \tilde{X}^T \tilde{X}$$

Uniqueness of the OLS solution

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (9)$$

- The OLS solution is available when the $D \times D$ square matrix $\mathbf{A} = \mathbf{X}^\top \mathbf{X}$ is invertible.
- Matrix inverse \mathbf{A}^{-1} can be computed by eigendecomposition:

$$\mathbf{A} = \mathbf{Q} \Lambda \mathbf{Q}^{-1} \implies \mathbf{A}^{-1} = \mathbf{Q}^{-1} \Lambda^{-1} \mathbf{Q}$$

where

- \mathbf{Q} is the square $D \times D$ matrix, whose i^{th} column is the i^{th} eigenvector
- \mathbf{Q} is also an **orthonormal** matrix, which means $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I} \implies \mathbf{Q}^\top = \mathbf{Q}^{-1}$
- Each eigenvector \mathbf{u} in \mathbf{Q} is orthonormal to itself: $\mathbf{u}^\top \mathbf{u} = 1$ and two different eigenvectors are orthogonal to each other: $\mathbf{u}^\top \mathbf{v} = 0$
- Λ is the diagonal matrix whose diagonal elements are the corresponding eigenvalues, i.e., $\Lambda_{i,i} = \lambda_i$
- Therefore, \mathbf{A} is not invertible if some of its eigenvalues are zeros, which can happen when two features are perfectly correlated, e.g., $\mathbf{x}_2 = 1 - \mathbf{x}_1$, meaning that the number of linearly independent columns (i.e., rank) is smaller than D .

* To make matrix invertible,
+ add noise ($+ \epsilon$)
→ invertible makes it more
stable for us to
compute w^*

$$\begin{bmatrix} \mathbf{x}^T \\ D \times N \end{bmatrix} \quad \begin{bmatrix} \quad \\ N \times D \end{bmatrix}$$

Time complexity

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\begin{bmatrix} \mathbf{x}^T \\ D \times N \end{bmatrix} \quad \begin{bmatrix} \quad \\ N \times C \end{bmatrix}$$

- The inner product $\mathbf{A} = \mathbf{X}^T \mathbf{X}$ takes $O(ND^2)$
- The inversion of the $D \times D$ matrix \mathbf{A}^{-1} takes $O(D^3)$
- Computing $\mathbf{X}^T \mathbf{y}$ takes $O(ND)$

Therefore, the total time complexity is $O(ND^2 + D^3)$.

This can be very expensive or infeasible for large D (e.g., 1 million SNPs or 20,000 genes) and/or large N (e.g., half million individuals). For this reason, although we can derive closed-form solution, *Stochastic Gradient Descent (SGD)* is used for large dataset (Module 4.4).

or Coordinate descent (CD)

Multivariate regression

We can also adapt the equation for multiple response variables. Instead of a response vector $\mathbf{y} \in \mathbb{R}^N$, we have a response matrix $\mathbf{Y} \in \mathbb{R}^{N \times C}$ for C response variables. The multivariate regression function is:

$$\mathbf{Y} = \mathbf{X}\mathbf{W} \quad (10)$$

$N \times D \times C$

where $\mathbf{W} \in \mathbb{R}^{D \times C}$.

The OLS solution for \mathbf{W} is then:

$$\mathbf{W} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \quad (11)$$

Note here the OLS coefficient \mathbf{w}_k for each response variable k is computed *independently* in the above solution. Therefore, the resulting OLS \mathbf{W} is identical to fitting each $D \times 1$ regression coefficient \mathbf{w}_k separately and then concatenate the C vectors together to form the matrix \mathbf{W} .

Outline

Objectives

Simple linear regression

Multiple linear regression

Probabilistic interpretation

Non-linear basis function

Summary

Probabilistic interpretation: Gaussian response variable

The general multivariate normal (MVN) distribution is:

$$p(\mathbf{y}|\mathbf{X}) = \mathcal{N}(\mathbf{y}|\mathbf{X}\mathbf{w}, \Sigma) = \det(2\pi\Sigma)^{-1/2} \exp\left(-\frac{1}{2} \underbrace{(\mathbf{y} - \mathbf{X}\mathbf{w})^\top}_{\mathbf{\epsilon}^\top} \Sigma^{-1} \underbrace{(\mathbf{y} - \mathbf{X}\mathbf{w})}_{\mathbf{\epsilon}}\right) \quad (12)$$

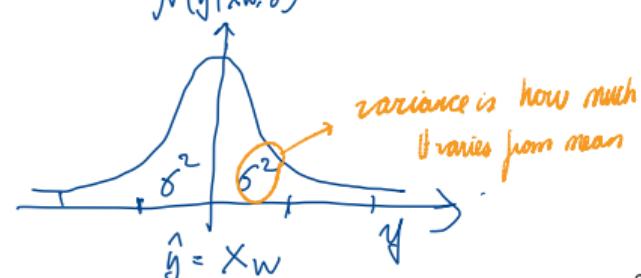
mean *variance* ↗ *quadratic form*

where Σ is a $N \times N$ covariance matrix between the N samples. If we assume the samples are *i.i.d.*, then Σ is a diagonal matrix $\sigma^2 \mathbf{I}$, where \mathbf{I} is an identity matrix:

$$\Sigma \stackrel{i.i.d.}{=} \begin{bmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}, \quad \Sigma^{-1} = \sigma^{-2} \mathbf{I}, \quad \det(2\pi\Sigma)^{-1/2} = (2\pi\sigma^2)^{-N/2}$$

independently identically distributed

where $\det(\mathbf{A})$ is the determinant of a square matrix.



The MVN can then be simplified as the product of individual Gaussians:

$$p(\mathbf{y}|\mathbf{X}) = (2\pi\sigma^2)^{-N/2} \exp\left(-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\mathbf{w})^\top(\mathbf{y} - \mathbf{X}\mathbf{w})\right)$$
$$\text{lab} = \ln c + \ln b = (2\pi\sigma^2)^{-N/2} \exp\left(-\frac{1}{2\sigma^2} \sum_n (y^{(n)} - \mathbf{x}^{(n)}\mathbf{w})^2\right)$$

Taking the logarithm of the likelihood, we have

$$\begin{aligned} \ln p(\mathbf{y}|\mathbf{X}) &= \underbrace{-\frac{N}{2} \ln(2\pi\sigma^2)}_{\text{constant w.r.t. } \mathbf{w}} - \sum_n \left(\frac{1}{2\sigma^2} (y^{(n)} - \mathbf{x}^{(n)}\mathbf{w})^2 \right) \\ &\propto -\frac{1}{2\sigma^2} \sum_n (y^{(n)} - \mathbf{x}^{(n)}\mathbf{w})^2 \\ &= \bigcirc \frac{1}{2\sigma^2} \boxed{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2} \end{aligned}$$

the smaller this is,
the higher the max

The last equation indicates that the log likelihood is proportional to the negative SSE.

Maximum likelihood estimation

Given that,

$$\ln p(\mathbf{y}|\mathbf{X}) \propto -\frac{1}{2\sigma^2} \sum_n (y^{(n)} - \mathbf{x}^{(n)} \mathbf{w})^2 = -\frac{1}{2\sigma^2} J(\mathbf{w})$$

Since σ^2 is a constant, minimizing SSE w.r.t. \mathbf{w} is equivalent to maximizing the Gaussian likelihood w.r.t. \mathbf{w} :

$$\arg \min_{\mathbf{w}} J(\mathbf{w}) = \arg \max_{\mathbf{w}} \ln p(\mathbf{y}|\mathbf{X})$$

The maximum likelihood estimator for \mathbf{w} is then identical to the OLS \mathbf{w} :

Optional:

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Bayesian inference of $P(\mathbf{w}|\mathbf{X}, \mathbf{y}) \propto \frac{P(\mathbf{y}|\mathbf{X}, \mathbf{w}) P(\mathbf{w})}{P(\mathbf{y}|\mathbf{X})} = \frac{P(\mathbf{y}|\mathbf{w}, \mathbf{X}) P(\mathbf{w})}{\int P(\mathbf{y}, \mathbf{w}|\mathbf{X}) d\mathbf{w}}$

Outline

Objectives

Simple linear regression

Multiple linear regression

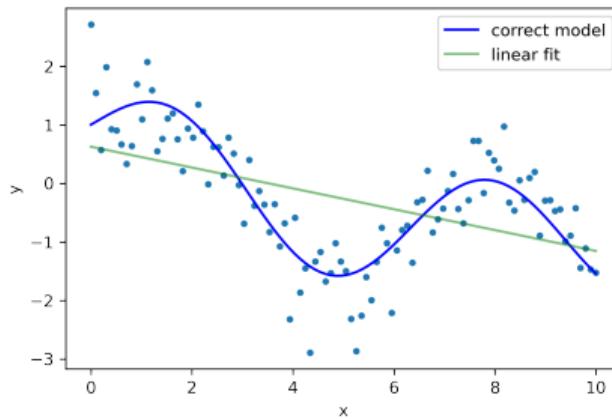
Probabilistic interpretation

Non-linear basis function

Summary

Fitting non-linear data by transforming the input features

Consider the toy dataset below. It is obvious that our attempt to model y as a linear function of $\hat{y} = w_0 + xw_1$ would produce a bad fit.



Idea: we can create more features using the given features. For example, we can use a M-degree polynomial function and treat each power degree as a standalone feature:

$$\hat{y} = \underbrace{x^0 w_0}_{\text{intercept}} + x^1 w_1 + x^2 w_2 + \dots + x^M w_M = \sum_{m=0}^M x^m w_m$$

Fitting non-linear data by transforming the input features

In general, we can transform input feature x with a (non-linear) **basis function** $\phi(x)$.
The multiple linear regression operates on the basis-transformed features:

$$\hat{y} = \sum_{m=0}^M \phi_m(x) w_m = \Phi(x)\mathbf{w} \quad (13)$$

$(x^{(M+1)})^{(M+1) \times 1}$

We then simply replace all of the occurrences of x with $\Phi(x)$ in the OLS solution:

$$\hat{\mathbf{w}} = (\Phi(x)^T \Phi(x))^{-1} \Phi(x)^T \mathbf{y}, \quad \text{where}$$

$$\Phi(x) = \begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) & \dots & \phi_M(x^{(1)}) \\ \phi_0(x^{(2)}) & \phi_1(x^{(2)}) & \dots & \phi_M(x^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^{(N)}) & \phi_1(x^{(N)}) & \dots & \phi_M(x^{(N)}) \end{bmatrix}$$

$N \times (M+1)$

Transforming high-dimensional features

- This can also be done on high-dimensional input feature $\mathbf{x}^{(n)}$ by transforming each feature with $x_d^{(n)}$ with say M -degree polynomial:

actual feature index d

$$\Phi(x_d^{(n)}) = [(x_d^{(n)})^0, \dots, (x_d^{(n)})^M]$$

Create M features out of d features

- We can then concatenate all transformed $D \times (M + 1)$ features

$$\mathbf{x}^{(n)} = [(x_1^{(n)})^0, \dots, (x_1^{(n)})^M, \dots, (x_D^{(n)})^0, \dots, (x_D^{(n)})^M] = [\Phi(x_1^{(n)}), \dots, \Phi(x_D^{(n)})]$$

($\mathbf{x} \in D \times M$)

- This creates a input matrix of dimension $N \times (D \times (M + 1))$:

$$\Phi(\mathbf{X}) = \begin{bmatrix} \Phi(x_1^{(1)}) & \Phi(x_2^{(1)}) & \dots & \Phi(x_D^{(1)}) \\ \Phi(x_1^{(2)}) & \Phi(x_2^{(2)}) & \dots & \Phi(x_D^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \Phi(x_1^{(N)}) & \Phi(x_2^{(N)}) & \dots & \Phi(x_D^{(N)}) \end{bmatrix}$$

$\overset{D \times M + 1}{\text{W}}$ $\rightarrow \hat{y}$ $(D \times M + 1) \times 1$ $N \times 1$

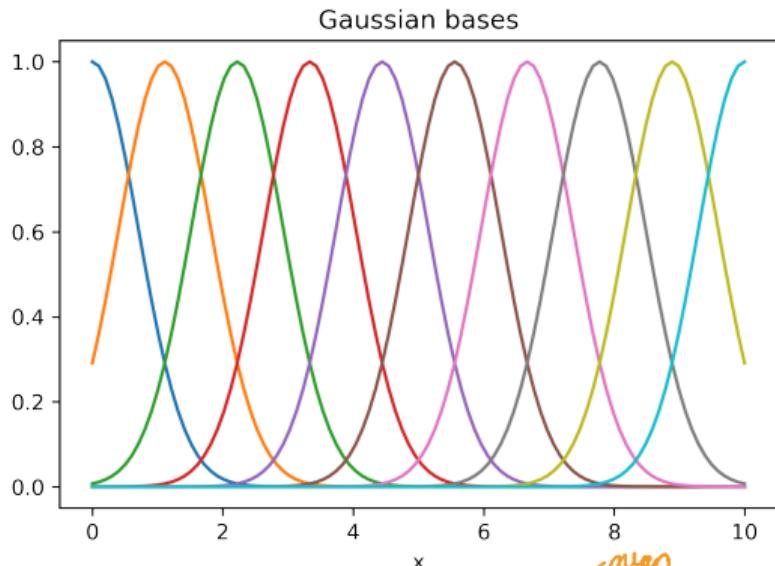
$\hat{y} = \mathbf{1} w_0 + \Phi(\mathbf{x}) W$

dimension?

redundant so $(x_d^{(n)})^0$ is not needed

Nonlinear basis functions

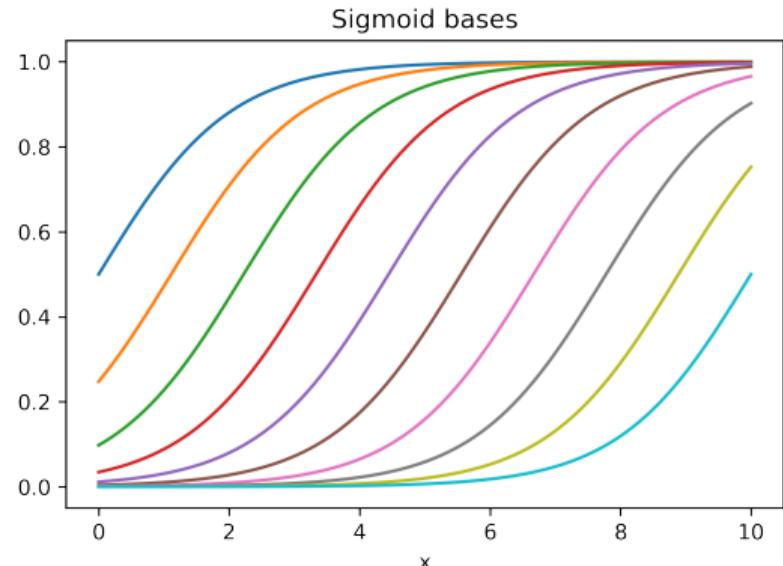
There are many nonlinear basis functions. Using scalar input $x \in \mathbb{R}$ as an example,



$$\phi_d(x) = \exp\left(-\frac{(x - \mu_d)^2}{2s^2}\right)$$

x

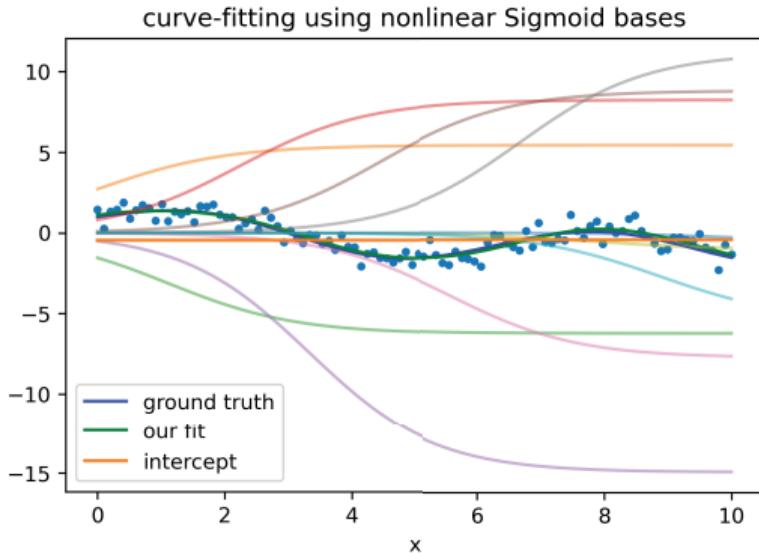
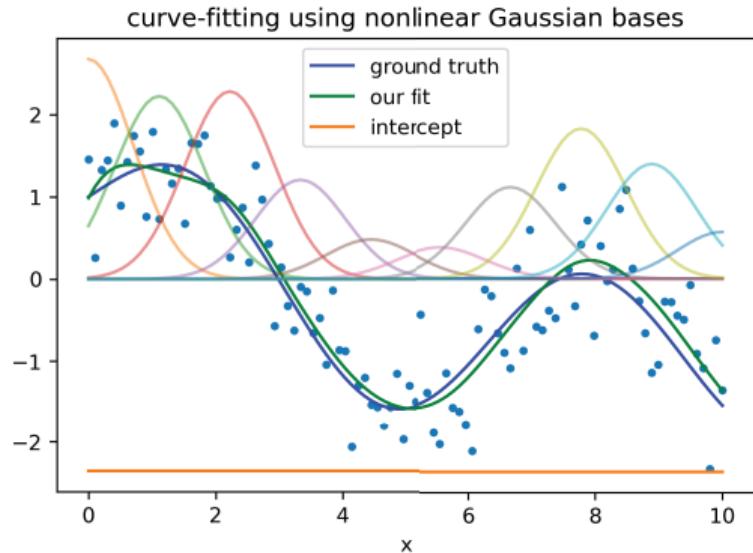
mean
variance



$$\phi_d(x) = \frac{1}{1 + \exp\left(-\frac{(x - \mu_d)^2}{s}\right)}$$

where each type of basis function has a different mean $\mu_d \in [0, 10]$ and $s = 1$.

Linear regression with nonlinear basis (See Colab for the implementations)



In both plots, the green curve (our fit) is the sum of weighted Gaussian bases (i.e., the colorful curves) plus the intercept:

* Creating more features, so that it fits better

$$\hat{y} = w_0 + \sum_d w_d \phi_d(\mathbf{x})$$

Summary

more than
1 feature

- Response variable y is modelled as a weighted linear sum of the features $\hat{y} = \mathbf{X}\mathbf{w}$
- We fit the model by minimizing the sum of squared errors (SSE)
 $\sum_n(y^{(n)} - \mathbf{X}^{(n)}\mathbf{w})^2 \rightarrow \text{achieve global optimal}$
- OLS solution: $\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ with $O(ND^2 + D^3)$ time complexity
- Minimizing SSE is equivalent to maximizing the log Gaussian likelihood given that the data points are *i.i.d.*
- Using non-linear basis functions to construct features can help model non-linear data. However, these non-linear basis functions are rigid and non-learnable.
- In Module 5.1, we will discuss neural networks which have *learnable* non-linear basis functions.