

Optimal MSE predictor is $\mathbb{E}(Y | X)$

Renato Assunção
ESRI and DCC/UFMG

A classical machine learning problem: regression

Regression example

Predict the price Y of an apartment from 30 numerical features such as area, X_2 age, X_3 number of rooms, and so on.

Let

$$\mathbf{X} = (X_1, X_2, \dots, X_{30}) \in \mathbb{R}^{30},$$

where X_1 is the area, X_2 is the age, X_3 is the number of rooms, etc.

Data structure

We have a sample of n observed data examples

That is, for the i -th apartment we observe:

$$(Y_i, X_{i1}, X_{i2}, \dots, X_{i,30}), \quad i = 1, \dots, n.$$

- Y_i : response variable (price of the i -th apartment).
- The remaining variables $x_{i1}, \dots, x_{i,30}$ are the 30 features of the i -th apartment.

Data structure

We have a sample of n observed data examples

That is, for the i -th apartment we observe:

$$(Y_i, X_{i1}, X_{i2}, \dots, X_{i,30}), \quad i = 1, \dots, n.$$

- Y_i : response variable (price of the i -th apartment).
- The remaining variables $x_{i1}, \dots, x_{i,30}$ are the 30 features of the i -th apartment.

Goal

Find a prediction rule $g(\mathbf{x})$ such that $Y - g(\mathbf{x})$ is small for future observations.

Seeking a predictor

We want a formula

$$g(\mathbf{x}) = g(x_1, x_2, \dots, x_{30})$$

that predicts the response Y well.

Seeking a predictor

We want a formula

$$g(\mathbf{x}) = g(x_1, x_2, \dots, x_{30})$$

that predicts the response Y well.

- For an observed unit, $Y - g(\mathbf{x})$ is the prediction error.
- What we really care about is prediction for a *new* case.
- This leads to a criterion to choose a good prediction rule g based on the expected prediction error.

Population prediction risk

A natural criterion is the mean squared prediction error:

$$\text{MSE}(g) = \mathbb{E}[(Y - g(\mathbf{X}))^2].$$

Population prediction risk

A natural criterion is the mean squared prediction error:

$$\text{MSE}(g) = \mathbb{E}[(Y - g(\mathbf{X}))^2].$$

- MSE = Mean Squared Error.
- We are looking at the expected value of the (squared) prediction error.
- One of the reasons to use the squared prediction error $(Y - g(\mathbf{X}))^2$ rather than its absolute value $|Y - g(\mathbf{X})|$ is that it is mathematically easier (more about that later).

Theoretical optimum

This set of slides aims to show that, among all infinite possible functions $g(\mathbf{X})$, the unique minimizer of

$$\mathbb{E}[(Y - g(\mathbf{X}))^2]$$

is to take the function g as

$$g(\mathbf{X}) = \mathbb{E}[Y \mid \mathbf{X}].$$

Theoretical optimum

This set of slides aims to show that, among all infinite possible functions $g(\mathbf{X})$, the unique minimizer of

$$\mathbb{E}[(Y - g(\mathbf{X}))^2]$$

is to take the function g as

$$g(\mathbf{X}) = \mathbb{E}[Y \mid \mathbf{X}].$$

Interpretation

The conditional mean is the best possible predictor under squared loss.

General Setup Motivation

- How to characterize a random variable Y ?
- Two "lists" ...
- Every random variable Y can be written as

$$Y = \mathbb{E}(Y) + (Y - \mathbb{E}(Y)) = \mathbb{E}(Y) + \varepsilon$$

- ε is the deviation of Y with respect to $\mathbb{E}(Y)$
- Is ε a random variable?
- How about $\mathbb{E}(Y)$?

General Setup Motivation

- Consider variables $(Y, \mathbf{X}) = (Y, X_1, X_2, \dots, X_k)$.
- The random vector (Y, \mathbf{X}) has a joint density $f(y, \mathbf{x})$.
- The joint density determines the conditional density $f(y|\mathbf{x})$.
- Goal: Predict the outcome Y using the features (variables) \mathbf{X} .
- We assume that we know the values of the random variables in \mathbf{X} .
- How to best predict Y given that we know that $\mathbf{X} = \mathbf{x}$.
- This requires knowing the conditional distribution $f(y|\mathbf{x})$ of the random variable $(Y|\mathbf{X} = \mathbf{x})$.

General Setup Motivation

Challenge

How to specify $f(y|\mathbf{x})$ for every possible configuration \mathbf{x} and every value y ?

Example: Apartment Price

$Y = \text{Price}$

$\mathbf{X} = (\text{Area } (X_1), \text{ #Quartos } (X_2), \text{ #Suítes } (X_3), \dots, \text{ Piscina? } (X_k), \dots)$

General Setup Motivation

Challenge

How to specify $f(y|\mathbf{x})$ for every possible configuration \mathbf{x} and every value y ?

Example: Apartment Price

Y = Price

\mathbf{X} = (Area (X_1), #Quartos (X_2), #Suítes (X_3), ..., Piscina? (X_k), ...)

For example, need the density $f(y|X_1 = 250m^2, X_2 = 3, \dots)$

Also need the density $f(Y|X_1 = 120m^2, X_2 = 2, \dots)$

...and for infinitely many other possibilities for \mathbf{X} .

Simplifying the Goal

Instead of finding the full distribution $f(y|\mathbf{x})$ for all \mathbf{x} , let's try something simpler:

- Predict Y using a single number $g(\mathbf{x})$.
- Given that \mathbf{X} is a certain specific configuration \mathbf{x} , we want a "formula", a function $g(\mathbf{x})$ that outputs a good prediction for the associated random Y .

Simplifying the Goal

Instead of finding the full distribution $f(y|\mathbf{x})$ for all \mathbf{x} , let's try something simpler:

- Predict Y using a single number $g(\mathbf{x})$.
- Given that \mathbf{X} is a certain specific configuration \mathbf{x} , we want a "formula", a function $g(\mathbf{x})$ that outputs a good prediction for the associated random Y .

- Let

$$g : \mathbb{R}^k \rightarrow \mathbb{R}$$
$$\mathbf{x} \mapsto g(\mathbf{x})$$

- We want $g(\mathbf{x}) \approx Y$ when $\mathbf{X} = \mathbf{x}$.

Simplifying the Goal

- For example, g could be a linear combination of the features in \mathbf{x} :

$$g(\mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

- Or it could be a more complex function based on \mathbf{x} such as, for example:

$$g(\mathbf{x}) = \beta_0 + \beta_1 \log(x_1) + \beta_2 \exp(-\beta_3 x_2 x_3) + x_4^{\beta_4}$$

Simplifying the Goal

- Very important: differentiate $g(\mathbf{x})$ from $g(\mathbf{X})$
- $g(\mathbf{x})$: a scalar associated to the realization \mathbf{x} (lowercase x)
- For example, if

$$\mathbf{X} = \mathbf{x}$$

with $\mathbf{x} = (250m^2, 3, \dots)$

Simplifying the Goal

- Very important: differentiate $g(\mathbf{x})$ from $g(\mathbf{X})$
- $g(\mathbf{x})$: a scalar associated to the realization \mathbf{x} (lowercase \mathbf{x})
- For example, if

$$\mathbf{X} = \mathbf{x}$$

with $\mathbf{x} = (250m^2, 3, \dots)$

Then, the predicted price may be

$$g(\mathbf{x}) = g(250m^2, 3, \dots) = 1.2M,$$

a scalar, a simple number.

Simplifying the Goal

- Very important: differentiate $g(\mathbf{x})$ from $g(\mathbf{X})$
- $g(\mathbf{x})$: a scalar associated to the realization \mathbf{x} (lowercase x)
- For example, if

$$\mathbf{X} = \mathbf{x}$$

with $\mathbf{x} = (250m^2, 3, \dots)$

Then, the predicted price may be

$$g(\mathbf{x}) = g(250m^2, 3, \dots) = 1.2M,$$

a scalar, a simple number.

- Now, we observe the event $\mathbf{X} = \mathbf{x}^*$ where $\mathbf{x}^* = (100m^2, 2, \dots) \neq \mathbf{x}$

Simplifying the Goal

- Very important: differentiate $g(\mathbf{x})$ from $g(\mathbf{X})$
- $g(\mathbf{x})$: a scalar associated to the realization \mathbf{x} (lowercase x)
- For example, if

$$\mathbf{X} = \mathbf{x}$$

with $\mathbf{x} = (250m^2, 3, \dots)$

Then, the predicted price may be

$$g(\mathbf{x}) = g(250m^2, 3, \dots) = 1.2M,$$

a scalar, a simple number.

- Now, we observe the event $\mathbf{X} = \mathbf{x}^*$ where $\mathbf{x}^* = (100m^2, 2, \dots) \neq \mathbf{x}$
- The predicted price may change to another scalar

$$g(\mathbf{x}^*) = g(100m^2, 2, \dots) = 0.8M$$

Simplifying the Goal

- How about $g(\mathbf{X})$, with a capital \mathbf{X} ?

Simplifying the Goal

- How about $g(\mathbf{X})$, with a capital \mathbf{X} ?
- This is not a scalar. It is a random variable because it is a mathematical transformation of the random vector \mathbf{X} .
- Being a random variable, $g(\mathbf{X})$ is composed of two things:

Simplifying the Goal

- How about $g(\mathbf{X})$, with a capital \mathbf{X} ?
- This is not a scalar. It is a random variable because it is a mathematical transformation of the random vector \mathbf{X} .
- Being a random variable, $g(\mathbf{X})$ is composed of two things:
- List of possible values and ...
- List of associated probabilities

Simplifying the Goal

- How about $g(\mathbf{X})$, with a capital \mathbf{X} ?
- This is not a scalar. It is a random variable because it is a mathematical transformation of the random vector \mathbf{X} .
- Being a random variable, $g(\mathbf{X})$ is composed of two things:
 - List of possible values and ...
 - List of associated probabilities
- Usually, the list of probabilities is impossible to know exactly unless \mathbf{X} has a very simple distribution, such as a multivariate Gaussian.

Simplifying the Goal

The Prediction Error

- Define the random prediction error: $\epsilon = Y - g(\mathbf{X})$
- Note that Y and $g(\mathbf{X})$ are both random variables.
- Goal: Choose g such that this random error ϵ is "small".
- For a fixed \mathbf{x} vector, $g(\mathbf{x})$ is a number, but Y is still random.
- So $\epsilon = Y - g(\mathbf{x})$ is a random variable, even when we have \mathbf{X} fixed at the observed value \mathbf{x}
- That is, even when we condition on the value of the random vector \mathbf{X} .

Measuring Error: MSE

How do we quantify "small" for the random error $\epsilon = Y - g(\mathbf{X})$?

- Rather than looking at $|\epsilon|$, look at the squared error:
 $\epsilon^2 = (Y - g(\mathbf{X}))^2$ (for math convenience).
- As ϵ^2 is random, it is hard to find a g such that this random error is always guaranteed small.

Measuring Error: MSE

How do we quantify "small" for the random error $\epsilon = Y - g(\mathbf{X})$?

- Rather than looking at $|\epsilon|$, look at the squared error:
 $\epsilon^2 = (Y - g(\mathbf{X}))^2$ (for math convenience).
- As ϵ^2 is random, it is hard to find a g such that this random error is always guaranteed small.
- We were able to solve another problem: rather than minimizing ϵ^2 100% of times, why don't we minimize the *expected* square error?
 $\mathbb{E}(\epsilon^2)$
- We may be happy to obtain a g such that, on average, the (square) error is small.
- Occasionally, we may have a large error but most of the times we will have a small error.

Measuring Error: MSE

Assume that (Y, \mathbf{X}) is a random vector.

In practice we will be able to know one part of this vector:

$$\mathbf{X} = \mathbf{x}$$

for some specific vector \mathbf{x} .

We will then use $g(\mathbf{x})$ to predict the unknown Y associated with this \mathbf{x} .

What should be g ?

Measuring Error: MSE

Optimization Problem

Find the function g that minimizes:

$$\text{MSE}(g) = \mathbb{E}[(Y - g(\mathbf{X}))^2]$$

The solution is $\hat{g} = \arg_g \min \mathbb{E}[(Y - g(\mathbf{X}))^2]$

Measuring Error: MSE

Optimization Problem

Find the function g that minimizes:

$$\text{MSE}(g) = \mathbb{E}[(Y - g(\mathbf{X}))^2]$$

The solution is $\hat{g} = \arg_g \min \mathbb{E}[(Y - g(\mathbf{X}))^2]$

The expectation has both random, Y and \mathbf{X} .

That is, MSE means

$$\mathbb{E}_{f(y,\mathbf{x})}[(Y - g(\mathbf{X}))^2] = \int \int_{\mathbb{R} \times \mathbb{R}^k} (y - g(\mathbf{x}))^2 f(y, \mathbf{x}) \, d\mathbf{x} dy$$

Solving the optimization problem

Using the iterated expectation:

$$\begin{aligned}\text{MSE}(g) &= \mathbb{E}_{Y, \mathbf{X}}[(Y - g(\mathbf{X}))^2] \\ &= \mathbb{E}_{\mathbf{X}} [\mathbb{E}_{Y|\mathbf{X}}[(Y - g(\mathbf{X}))^2 | \mathbf{X}]]\end{aligned}$$

Solving the optimization problem

Using the iterated expectation:

$$\begin{aligned}\text{MSE}(g) &= \mathbb{E}_{Y, \mathbf{X}}[(Y - g(\mathbf{X}))^2] \\ &= \mathbb{E}_{\mathbf{X}} [E_{Y|\mathbf{X}}[(Y - g(\mathbf{X}))^2|\mathbf{X}]] \\ &= \int_{\mathbb{R}^k} \underbrace{\left[\int_{-\infty}^{\infty} (y - g(\mathbf{x}))^2 f(y|\mathbf{x}) dy \right]}_{h(\mathbf{x}) = E_{Y|\mathbf{X}}[(Y - g(\mathbf{x}))^2|\mathbf{X}=\mathbf{x}]} f(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Solving the optimization problem

Using the iterated expectation:

$$\begin{aligned}\text{MSE}(g) &= \mathbb{E}_{Y, \mathbf{X}}[(Y - g(\mathbf{X}))^2] \\ &= \mathbb{E}_{\mathbf{X}} [E_{Y|\mathbf{X}}[(Y - g(\mathbf{X}))^2|\mathbf{X}]] \\ &= \int_{\mathbb{R}^k} \underbrace{\left[\int_{-\infty}^{\infty} (y - g(\mathbf{x}))^2 f(y|\mathbf{x}) dy \right]}_{h(\mathbf{x}) = E_{Y|\mathbf{X}}[(Y - g(\mathbf{x}))^2|\mathbf{X}=\mathbf{x}]} f(\mathbf{x}) d\mathbf{x}\end{aligned}$$

$h(\mathbf{x}) \geq 0$ for all fixed value of \mathbf{x} because the integrand $(y - g(\mathbf{x}))^2 f(y|\mathbf{x}) \geq 0$ for all y and \mathbf{x} .

Solving the optimization problem

Using the iterated expectation:

$$\begin{aligned}\text{MSE}(g) &= \mathbb{E}_{Y, \mathbf{X}}[(Y - g(\mathbf{X}))^2] \\ &= \mathbb{E}_{\mathbf{X}} \left[E_{Y|\mathbf{X}}[(Y - g(\mathbf{X}))^2 | \mathbf{X}] \right] \\ &= \int_{\mathbb{R}^k} \underbrace{\left[\int_{-\infty}^{\infty} (y - g(\mathbf{x}))^2 f(y|\mathbf{x}) dy \right]}_{h(\mathbf{x}) = E_{Y|\mathbf{X}}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]} f(\mathbf{x}) d\mathbf{x}\end{aligned}$$

$h(\mathbf{x}) \geq 0$ for all fixed value of \mathbf{x} because the integrand $(y - g(\mathbf{x}))^2 f(y|\mathbf{x}) \geq 0$ for all y and \mathbf{x} .

Hence, we can minimize $\text{MSE}(g)$ (the double integral) by minimizing $h(\mathbf{x})$ at each possible value of \mathbf{x} . Integrating the smallest possible $h(\mathbf{x})$ guarantees that the double integral is minimum.

Optimal Predictor: Conditional Expectation

We want to minimize the inner expectation for each fixed \mathbf{x} :

$$\min_{g(\mathbf{x})} h(\mathbf{x}) = \min_{g(\mathbf{x})} \mathbb{E}_{Y|\mathbf{X}}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]$$

Optimal Predictor: Conditional Expectation

We want to minimize the inner expectation for each fixed \mathbf{x} :

$$\min_{g(\mathbf{x})} h(\mathbf{x}) = \min_{g(\mathbf{x})} \mathbb{E}_{Y|\mathbf{X}}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]$$

Let $\mu(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$.

We will show now that the optimal $g(\mathbf{x})$ must be $\mu(\mathbf{x})$

Start with the old trick that seems not to change anything, but that makes all the difference: add and subtract $\mu(\mathbf{x})$.

$$h(\mathbf{x}) = \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] = \mathbb{E}[(Y - \mu(\mathbf{x}) + \mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]$$

Proof:

$$\begin{aligned}h(\mathbf{x}) &= \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\ &= \mathbb{E}[(Y - \mu(\mathbf{x}) + \mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]\end{aligned}$$

Proof:

$$\begin{aligned}h(\mathbf{x}) &= \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}) + \mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&\quad + 2\mathbb{E}[(Y - \mu(\mathbf{x}))(\mu(\mathbf{x}) - g(\mathbf{x})) | \mathbf{X} = \mathbf{x}] \\&\quad + \mathbb{E}[(\mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}]\end{aligned}$$

Proof:

$$\begin{aligned}h(\mathbf{x}) &= \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}) + \mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&\quad + 2\mathbb{E}[(Y - \mu(\mathbf{x}))(\mu(\mathbf{x}) - g(\mathbf{x})) | \mathbf{X} = \mathbf{x}] \\&\quad + \mathbb{E}[(\mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) \\&\quad + 2(\mu(\mathbf{x}) - g(\mathbf{x})) \underbrace{\mathbb{E}[Y - \mu(\mathbf{x}) | \mathbf{X} = \mathbf{x}]}_{=0} \\&\quad + (\mu(\mathbf{x}) - g(\mathbf{x}))^2\end{aligned}$$

Proof:

$$\begin{aligned}h(\mathbf{x}) &= \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}) + \mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{E}[(Y - \mu(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&\quad + 2\mathbb{E}[(Y - \mu(\mathbf{x}))(\mu(\mathbf{x}) - g(\mathbf{x})) | \mathbf{X} = \mathbf{x}] \\&\quad + \mathbb{E}[(\mu(\mathbf{x}) - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] \\&= \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) \\&\quad + 2(\mu(\mathbf{x}) - g(\mathbf{x})) \underbrace{\mathbb{E}[Y - \mu(\mathbf{x}) | \mathbf{X} = \mathbf{x}]}_{=0} \\&\quad + (\mu(\mathbf{x}) - g(\mathbf{x}))^2 \\&= \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) + (\mu(\mathbf{x}) - g(\mathbf{x}))^2\end{aligned}$$

Conclusion:

We found that

$$h(\mathbf{x}) = \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] = \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) + (\mu(\mathbf{x}) - g(\mathbf{x}))^2$$

Conclusion:

We found that

$$h(\mathbf{x}) = \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] = \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) + (\mu(\mathbf{x}) - g(\mathbf{x}))^2$$

The first term does not depend on $g(\mathbf{x})$. Only the second term involves $g(\mathbf{x})$.

Hence, this expression is minimized when $(\mu(\mathbf{x}) - g(\mathbf{x}))^2 = 0$.

This means that we must have $g(\mathbf{x}) = \mu(\mathbf{x})$.

Conclusion:

We found that

$$h(\mathbf{x}) = \mathbb{E}[(Y - g(\mathbf{x}))^2 | \mathbf{X} = \mathbf{x}] = \mathbb{V}(Y | \mathbf{X} = \mathbf{x}) + (\mu(\mathbf{x}) - g(\mathbf{x}))^2$$

The first term does not depend on $g(\mathbf{x})$. Only the second term involves $g(\mathbf{x})$.

Hence, this expression is minimized when $(\mu(\mathbf{x}) - g(\mathbf{x}))^2 = 0$.

This means that we must have $g(\mathbf{x}) = \mu(\mathbf{x})$.

Conclusion

The function $g(\mathbf{x})$ that minimizes the MSE is the conditional expectation:

$$g(\mathbf{x}) = \mu(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$$

The Solution and Practical Challenges

Transforming from the scalar $g(\mathbf{x}) \Rightarrow$ to the random variable $g(\mathbf{X})$.

Optimal MSE Predictor

The best prediction random function $g(\mathbf{X})$ in terms of minimizing Mean Squared Error is the conditional expectation function:

$$g(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

The Solution and Practical Challenges

Transforming from the scalar $g(\mathbf{x}) \Rightarrow$ to the random variable $g(\mathbf{X})$.

Optimal MSE Predictor

The best prediction random function $g(\mathbf{X})$ in terms of minimizing Mean Squared Error is the conditional expectation function:

$$g(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

Think of the random variable $g(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$ as a prediction rule: when we receive the information that the random vector (two things...) \mathbf{X} has been instantiated as the specific vector \mathbf{x} , we predict the random unknown Y using $\hat{Y} = g(\mathbf{x})$.

The Solution and Practical Challenges

Transforming from the scalar $g(\mathbf{x}) \Rightarrow$ to the random variable $g(\mathbf{X})$.

Optimal MSE Predictor

The best prediction random function $g(\mathbf{X})$ in terms of minimizing Mean Squared Error is the conditional expectation function:

$$g(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$$

Think of the random variable $g(\mathbf{X}) = \mathbb{E}[Y|\mathbf{X}]$ as a prediction rule: when we receive the information that the random vector (two things...) \mathbf{X} has been instantiated as the specific vector \mathbf{x} , we predict the random unknown Y using $\hat{Y} = g(\mathbf{x})$.

Fantastic news: when aiming at minimizing the prediction error of Y , our objective should always be to find a good approximation to $E[Y|\mathbf{X}]$. The issue is that this approximation can be hard to obtain.

The Solution and Practical Challenges

Practical Difficulty: Estimating $\mu(\mathbf{X})$

Imagine we want to estimate $\mathbb{E}[Y | \text{Area} = 250, \text{Quartos} = 3, \dots]$.

That is, we want an approximation to $\mu(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$ where $\mathbf{x} = (250, 3, \dots)$

With a large dataset with many observed y 's and \mathbf{x} 's we might:

- Find all observations with these features: Area=250, Quartos=3, etc
- Average their corresponding Y values.
- This empirical average is approx $\mathbb{E}[Y | \text{Area} = 250, \text{Quartos} = 3, \dots]$.

The Solution and Practical Challenges

Practical Difficulty: Estimating $\mu(\mathbf{X})$

Imagine we want to estimate $\mathbb{E}[Y | \text{Area} = 250, \text{Quartos} = 3, \dots]$.

That is, we want an approximation to $\mu(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$ where $\mathbf{x} = (250, 3, \dots)$

With a large dataset with many observed y 's and \mathbf{x} 's we might:

- Find all observations with these features: Area=250, Quartos=3, etc
- Average their corresponding Y values.
- This empirical average is approx $\mathbb{E}[Y | \text{Area} = 250, \text{Quartos} = 3, \dots]$.

Problem:

- With many features, it is hard to find a large number of matches, even allowing for a certain approximation (such as, rather than area = 250, we may take area between 220 and 280 squared meters).
- Even having all features as discrete, specific combinations might be rare or non-existent in the sample ("curse of dimensionality").

Approximating $\mu(\mathbf{x})$

We need to find a better way to *approximate* $E[Y|\mathbf{X} = \mathbf{x}]$.

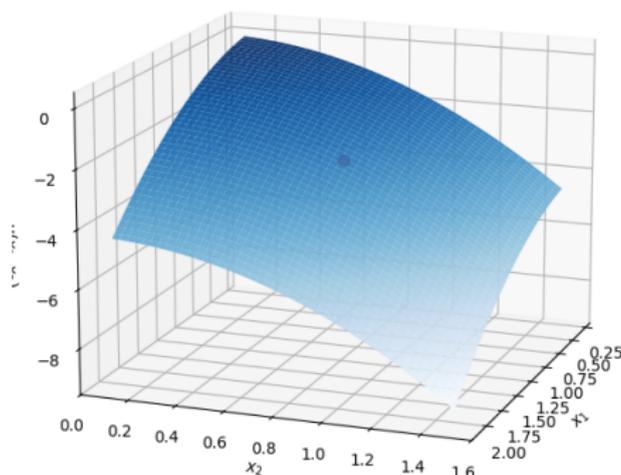
Approximating $\mu(\mathbf{x})$

We need to find a better way to *approximate* $E[Y|\mathbf{X} = \mathbf{x}]$.

Make an assumption:

$\mu(\mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ is a smooth function of $\mathbf{x} = (x_1, x_2, \dots, x_k)$.

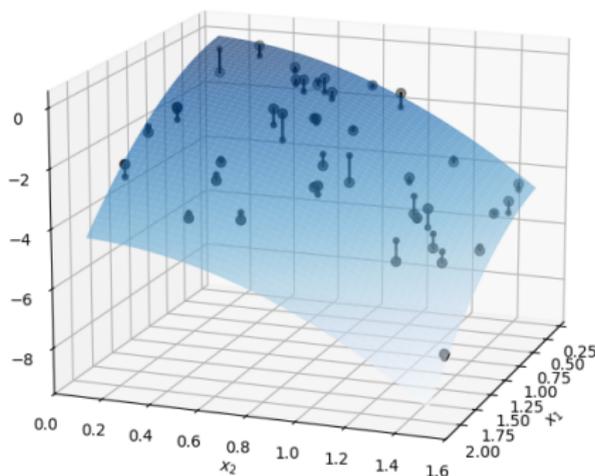
True surface $\mu(x_1, x_2)$



Approximating $\mu(\mathbf{x})$

The observed random Y spread around the surface $\mu(\mathbf{x})$:

Observed data: $Y_i = \mu(x_{1i}, x_{2i}) + \varepsilon_i$



Approximating $\mu(\mathbf{x})$

We can always write:

$$(Y|\mathbf{X} = \mathbf{x}) = \mathbb{E}[Y|\mathbf{X} = \mathbf{x}] + \epsilon = \mu(\mathbf{x}) + \epsilon$$

where ϵ is the random deviation from the conditional mean $\mu(\mathbf{x})$.

Components

- $\mu(\mathbf{x})$: Deterministic part, not random, depends smoothly on \mathbf{x} . Captures the underlying relationship.
- ϵ : Random part, represents noise or unmodeled factors that make the observation deviate from its expected value $\mu(\mathbf{x})$.

Approximating $\mu(\mathbf{x})$

Goal: Find a good, tractable mathematical approximation for the potentially complex function $\mu(\mathbf{x})$.

Approximating $\mu(\mathbf{x})$

Goal: Find a good, tractable mathematical approximation for the potentially complex function $\mu(\mathbf{x})$.

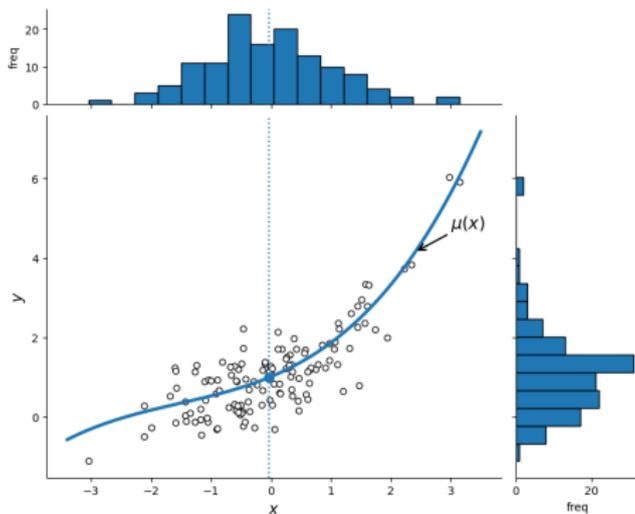
Polynomials are the simplest math functions: very easy to take derivatives and to integrate.

Mathematics has a method to approximate ANY complex math function $\mu(\mathbf{x})$ using polynomials: Taylor expansion.

This will serve as an initial justification to adopt the linear regression model.

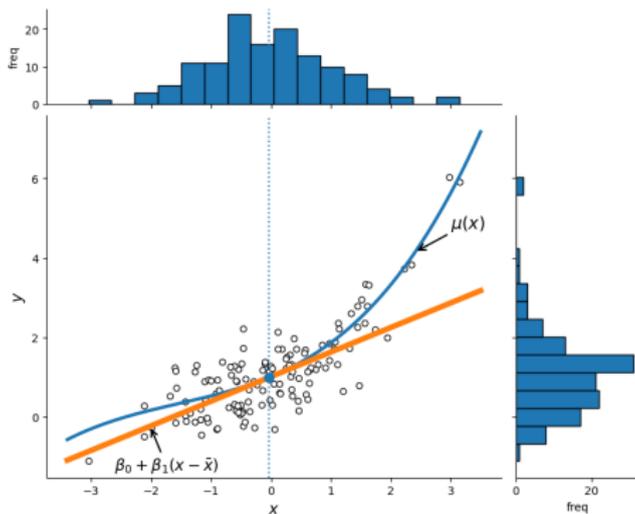
Taylor Approximation: one single feature x

Sample from (X, Y) and conditional mean curve $\mu(x) = E(Y | X = x)$



Taylor Approximation: one single feature x

Conditional mean curve and first-order Taylor approximation



Taylor Approximation: one single feature x

First-order Taylor expansion of $\mu(x)$: tangent line.

- Select a value of x around which we will make the approximation (expansion): the mean \bar{x} .
- Get the tangent line passing through $(\bar{x}, \mu(\bar{x}))$

$$\begin{aligned} Y &= \mu(x) + \varepsilon \\ &\approx \mu(\bar{x}) + \left[\frac{\partial \mu}{\partial x} \Big|_{x=\bar{x}} \right] (x - \bar{x}) + \varepsilon \\ &= \beta_0 + \beta_1(x - \bar{x}) + \varepsilon \end{aligned}$$

Taylor Approximation: one single feature x

$$\mu(x) = 1.0 + 0.75x + 0.22(x - 0.3)^2 + 0.04(x - 0.3)^3,$$

with $x \in (-3, 3)$ and $\bar{x} = 0.0$,

then: $\mu(\bar{x}) = \mu(0.0) \approx 1.02$ and

$$\begin{aligned} \left. \frac{\partial \mu}{\partial x} \right|_{x=0} &= \mu'(x) \Big|_{x=0} \\ &= (0.75 + 0.44(x - 0.3) + 0.12(x - 0.3)^2) \Big|_{x=0} \approx 0.63 \end{aligned}$$

Hence, the tangent line is $\hat{\mu}(x) = 1.02 + 0.63 * (x - 0.0)$

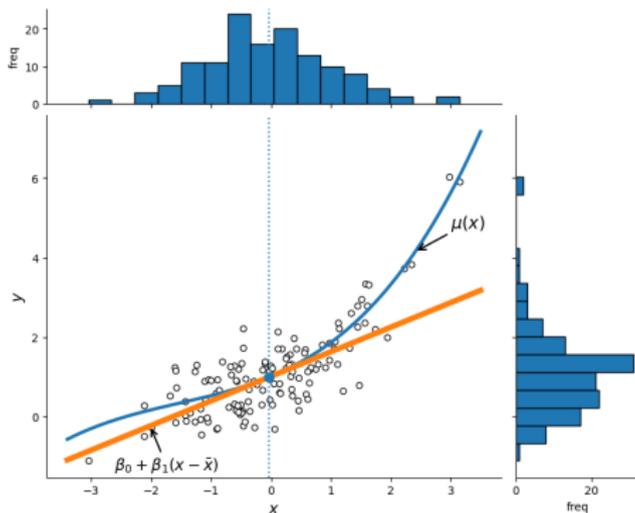
Taylor Approximation: one single feature x

Suppose

$$\mu(x) = 1.0 + 0.75x + 0.22(x - 0.3)^2 + 0.04(x - 0.3)^3$$

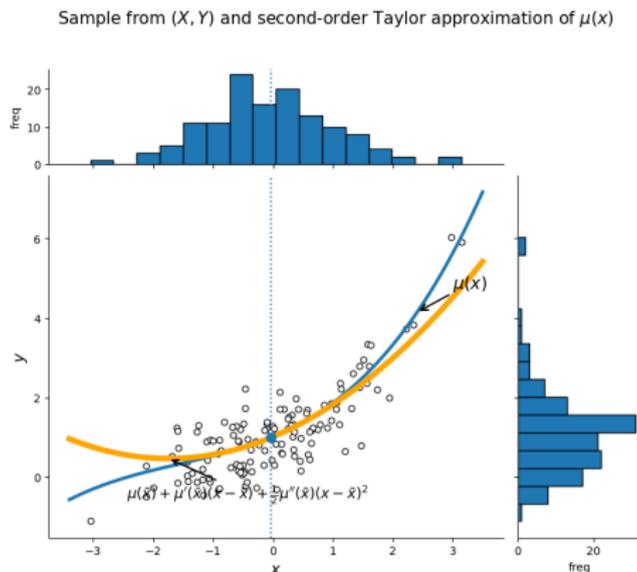
The tangent line is $\hat{\mu}(x) = 1.02 + 0.63 * (x - 0.0)$

Conditional mean curve and first-order Taylor approximation



Second-order Taylor Approximation

Second-order Taylor expansion of $\mu(x)$: tangent parabola around $\bar{x} \approx 0.0$.



Second-order Taylor Approximation

$$Y = \mu(x) + \varepsilon$$

$$\approx \mu(\bar{x}) + \left[\frac{\partial \mu}{\partial x} \Big|_{x=\bar{x}} \right] (x - \bar{x}) + \frac{1}{2} \left[\frac{\partial^2 \mu}{\partial x^2} \Big|_{x=\bar{x}} \right] (x - \bar{x})^2 + \varepsilon$$

$$= \beta_0 + \beta_1(x - \bar{x}) + \beta_2(x - \bar{x})^2 + \varepsilon$$

Taylor Approximation: one single feature x

$$\mu(x) = 1.0 + 0.75x + 0.22(x - 0.3)^2 + 0.04(x - 0.3)^3$$

with $\bar{x} = 0.0$, and $\mu'(0) = 0.63$

$$\begin{aligned}\left. \frac{\partial^2 \mu}{\partial x^2} \right|_{x=0} &= \mu''(x) \Big|_{x=0} \\ &= (0.44 + 0.24(x - 0.3)^2) \Big|_{x=0} \approx 0.37\end{aligned}$$

Hence, the tangent parabola is

$$\hat{\mu}(x) = 1.02 + 0.63 * (x - 0.0) + \frac{1}{2}0.37(x - 0.0)^2$$

Note a nice property: we did not change the linear component when creating the tangent parabola. We simply "added " a 2nd degree component.

Optimality of the Taylor approximation

In a sense, Taylor polynomial is the best **local** polynomial approximation you can get for a function $\mu(x)$.

Avoiding technicalities, the idea is that the parabola (2nd degree polynomial) that is the closest to $\mu(x)$ around a generic point a is the second-degree Taylor expansion around a .

This optimality is valid for high-degree polynomials: Taylor expansion of degree k is the best approximating k -th degree polynomial around a .

For details, see an Analysis book such as Elon Lages Lima.

Taylor approximation for multivariate $g(\mathbf{x})$

Usually, $\mu(\mathbf{x})$ is a function of many (k) features or covariates.

We have

$$\begin{aligned}\mu &: \mathbb{R}^k \rightarrow \mathbb{R} \\ \mathbf{x} &\mapsto \mu(\mathbf{x})\end{aligned}$$

How is the Taylor expansion in this case?

Consider the multi-dimensional tangent plane first, equivalent to the tangent line.

The gradient vector

The equivalent concept of the simple univariate derivative $\mu'(x)$ for a smooth k -dimensional function $g(\mathbf{x})$ is the gradient vector composed of the partial derivatives:

$$\nabla\mu(\mathbf{x}) = \begin{pmatrix} \frac{\partial\mu(\mathbf{x})}{\partial x_1} \\ \frac{\partial\mu(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial\mu(\mathbf{x})}{\partial x_k} \end{pmatrix}$$

Vectors will always be **column vectors** in this course. If you need a row vector, take the transpose:

$$\nabla\mu(\mathbf{x})^\top = \left(\frac{\partial\mu(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial\mu(\mathbf{x})}{\partial x_k} \right)$$

The gradient vector

The gradient $\nabla\mu(\mathbf{x})$ is a function of \mathbf{x} .

It changes depending on which point \mathbf{x} we evaluated it.

We denote the evaluation at the point $\mathbf{p} = (p_1, \dots, p_k)$ by:

$$\nabla\mu(\mathbf{x})|_{\mathbf{p}} = \left(\begin{array}{c} \frac{\partial\mu(\mathbf{x})}{\partial x_1} \\ \frac{\partial\mu(\mathbf{x})}{\partial x_2} \\ \vdots \\ \frac{\partial\mu(\mathbf{x})}{\partial x_k} \end{array} \right) \Big|_{\mathbf{p}}$$

First-order Taylor expansion: Tangent plane

A standard way to approximate a smooth function $g(\mathbf{x})$ locally around a point \mathbf{p} is using its Taylor expansion. Let $g : \mathbb{R}^k \rightarrow \mathbb{R}$.

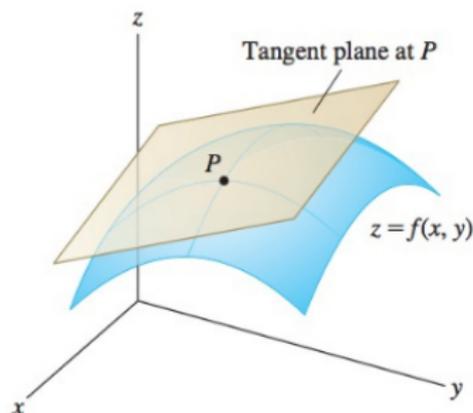
First-order Taylor expansion: Tangent plane

A standard way to approximate a smooth function $g(\mathbf{x})$ locally around a point \mathbf{p} is using its Taylor expansion. Let $g : \mathbb{R}^k \rightarrow \mathbb{R}$.

Taylor Expansion (First Order)

$$g(\mathbf{x}) \approx g(\mathbf{p}) + (\nabla g|_{\mathbf{p}})^T (\mathbf{x} - \mathbf{p})$$

where $\nabla g|_{\mathbf{p}}$ is the gradient vector of g evaluated at \mathbf{p} .



Example

$$g(\mathbf{x}) = e^{-3x_1^2 + 2x_1x_2 + x_2^2} \text{ around } \mathbf{p} = (1, 1)$$

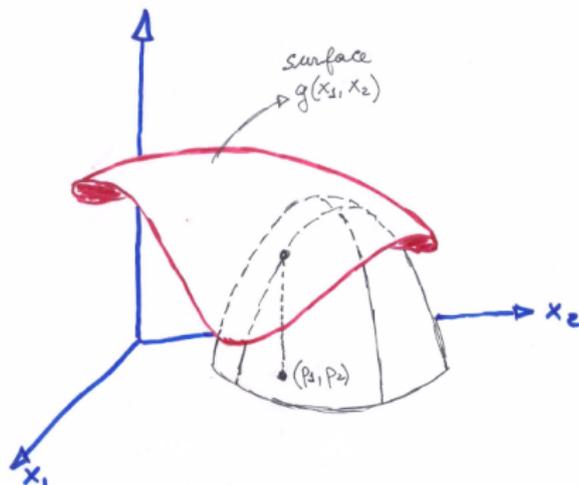
$$g(1, 1) = 1.$$

$$\nabla g(\mathbf{x})|_{(1,1)} = \begin{pmatrix} \frac{\partial g(\mathbf{x})}{\partial x_1} \\ \frac{\partial g(\mathbf{x})}{\partial x_2} \end{pmatrix} \Big|_{\mathbf{p}} = \begin{pmatrix} g(\mathbf{x})(-6x_1 + 2x_2) \\ g(\mathbf{x})(2x_1 + 2x_2) \end{pmatrix} \Big|_{\mathbf{p}} = \begin{pmatrix} -4 \\ 4 \end{pmatrix}$$

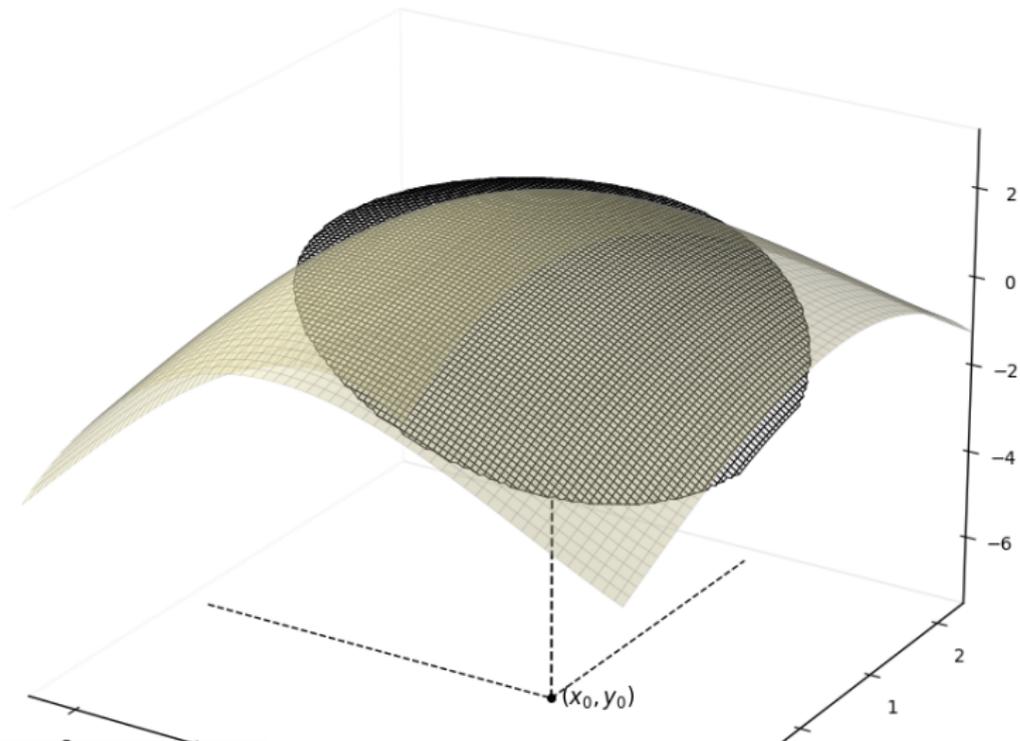
$$g(\mathbf{x}) \approx 1 + \begin{pmatrix} -4 & 4 \end{pmatrix} \begin{pmatrix} x_1 - 1 \\ x_2 - 1 \end{pmatrix} = 1 - 4x_1 + 4x_2$$

Second-order Taylor expansion: Tangent paraboloid

Rather than a simple plane, we can approximate the surface $g(\mathbf{x})$ locally around a point \mathbf{p} by the best possible paraboloid using its Taylor expansion.



Second-order Taylor expansion: Tangent paraboloid



Second-order Taylor expansion: Tangent paraboloid

Taylor Expansion (Second Order)

$$g(\mathbf{x}) \approx g(\mathbf{p}) + (\nabla g(\mathbf{x})|_{\mathbf{p}})^T (\mathbf{x} - \mathbf{p}) + \frac{1}{2} (\mathbf{x} - \mathbf{p})^T (H(\mathbf{x})|_{\mathbf{p}}) (\mathbf{x} - \mathbf{p})$$

where $H(\mathbf{x})|_{\mathbf{p}}$ is the Hessian matrix (matrix of second partial derivatives) of g evaluated at \mathbf{p} .

$$H(\mathbf{x}) = \begin{pmatrix} \frac{\partial^2 g(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_2} \\ \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_2} & \frac{\partial^2 g(\mathbf{x})}{\partial x_2^2} \end{pmatrix}$$

The final approximation expression is a polynomial in x_1 and x_2 involving up to the second-degree powers and the cross product $x_1 x_2$.

Second-order Taylor expansion: Tangent paraboloid

$$g(\mathbf{x}) \approx g(\mathbf{p}) + (\nabla g(\mathbf{x})|_{\mathbf{p}})^{\top}(\mathbf{x} - \mathbf{p}) + \frac{1}{2}(\mathbf{x} - \mathbf{p})^{\top}(H(\mathbf{x})|_{\mathbf{p}})(\mathbf{x} - \mathbf{p})$$

Let:

- $(\nabla g(\mathbf{x})|_{\mathbf{p}})^{\top} = (d1, d2)^{\top}$
- $H(\mathbf{x})|_{\mathbf{p}} = \left(\begin{array}{cc} \frac{\partial^2 g(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_2} \\ \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_2} & \frac{\partial^2 g(\mathbf{x})}{\partial x_2^2} \end{array} \right) \Big|_{\mathbf{p}} = \begin{pmatrix} a & b \\ b & c \end{pmatrix}$

Then (verify this)

$$g(\mathbf{x}) \approx g(\mathbf{p}) + d1(x_1 - p_1) + d2(x_2 - p_2) + \frac{1}{2}(a(x_1 - p_1)^2 + 2b(x_1 - p_1)(x_2 - p_2) + c(x_2 - p_2)^2)$$

Example

$$g(\mathbf{x}) = e^{-3x_1^2 + 2x_1x_2 + x_2^2} \text{ around } \mathbf{p} = (1, 1)$$

$$g(1, 1) = 1.$$

$$\nabla g(\mathbf{x})|_{(1,1)} = (-4, 4)^\top$$

$$H(\mathbf{x})|_{(1,1)=g(\mathbf{x})} \begin{pmatrix} -6 + (-6x_1 + 2x_2)^2 & -10x_1 + 6x_2 \\ -10x_1 + 6x_2 & 2 + 4(x_1 + x_2)^2 \end{pmatrix} \Big|_{(1,1)} = \begin{pmatrix} 10 & -4 \\ -4 & 18 \end{pmatrix}$$

$$g(\mathbf{x}) \approx 1 + (-4 \quad 4) \begin{pmatrix} x_1 - 1 \\ x_2 - 1 \end{pmatrix} + (x_1 - 1, x_2 - 1) \begin{pmatrix} 10 & -4 \\ -4 & 18 \end{pmatrix} \begin{pmatrix} x_1 - 1 \\ x_2 - 1 \end{pmatrix}$$

$$g(\mathbf{x}) \approx 1 - 4(x_1 - x_2) + 10(x_1 - 1)^2 + -8(x_1 - 1)(x_2 - 1) + 18(x_2 - 1)^2$$

Applying Taylor Expansion to $\mu(\mathbf{x})$

Let's approximate the conditional expectation function

$\mu(\mathbf{x}) = E[Y|\mathbf{X} = \mathbf{x}]$ using a first-order Taylor expansion around a central point, often the empirical mean $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_k)$.

Applying Taylor Expansion to $\mu(\mathbf{x})$

Let's approximate the conditional expectation function $\mu(\mathbf{x}) = E[Y|\mathbf{X} = \mathbf{x}]$ using a first-order Taylor expansion around a central point, often the empirical mean $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_k)$.

First-Order Approximation around $\bar{\mathbf{X}}$

$$\mu(\mathbf{x}) \approx \mu(\bar{\mathbf{x}}) + (\nabla\mu|_{\bar{\mathbf{x}}})^T (\mathbf{x} - \bar{\mathbf{x}})$$

Applying Taylor Expansion to $\mu(\mathbf{x})$

Let's approximate the conditional expectation function $\mu(\mathbf{x}) = E[Y|\mathbf{X} = \mathbf{x}]$ using a first-order Taylor expansion around a central point, often the empirical mean $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_k)$.

First-Order Approximation around $\bar{\mathbf{X}}$

$$\mu(\mathbf{x}) \approx \mu(\bar{\mathbf{x}}) + (\nabla \mu|_{\bar{\mathbf{x}}})^T (\mathbf{x} - \bar{\mathbf{x}})$$

Expanding this:

$$\mu(\mathbf{x}) \approx \mu(\bar{\mathbf{x}}) + \sum_{j=1}^k \left[\left. \frac{\partial \mu}{\partial x_j} \right|_{\bar{\mathbf{x}}} \right] (x_j - \bar{x}_j)$$

Applying Taylor Expansion to $\mu(\mathbf{x})$

Let's approximate the conditional expectation function $\mu(\mathbf{x}) = E[Y|\mathbf{X} = \mathbf{x}]$ using a first-order Taylor expansion around a central point, often the empirical mean $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_k)$.

First-Order Approximation around $\bar{\mathbf{X}}$

$$\mu(\mathbf{x}) \approx \mu(\bar{\mathbf{x}}) + (\nabla\mu|_{\bar{\mathbf{x}}})^T (\mathbf{x} - \bar{\mathbf{x}})$$

Expanding this:

$$\mu(\mathbf{x}) \approx \mu(\bar{\mathbf{x}}) + \sum_{j=1}^k \left[\frac{\partial\mu}{\partial x_j} \Big|_{\bar{\mathbf{x}}} \right] (x_j - \bar{x}_j)$$

Rearranging terms:

$$\mu(\mathbf{x}) \approx \underbrace{\left(\mu(\bar{\mathbf{x}}) - \sum_{j=1}^k \left[\frac{\partial\mu}{\partial x_j} \Big|_{\bar{\mathbf{x}}} \right] \bar{x}_j \right)}_{\beta_0 \text{ (Intercept)}} + \sum_{j=1}^k \underbrace{\left[\frac{\partial\mu}{\partial x_j} \Big|_{\bar{\mathbf{x}}} \right]}_{\beta_j \text{ (Slope for } X_j)} x_j$$

Applying Taylor Expansion to $\mu(\mathbf{x})$

Result

The first-order Taylor approximation of $E[Y|\mathbf{X} = \mathbf{x}]$ around $\bar{\mathbf{x}}$ has the form of a **linear function**:

$$\mu(\mathbf{x}) \approx \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k$$

The Linear Model

The first-order Taylor approximation motivates the linear model.

$$E[Y|\mathbf{X} = \mathbf{x}] = \mu(\mathbf{x}) \approx \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

The Linear Model

The first-order Taylor approximation motivates the linear model.

$$E[Y|\mathbf{X} = \mathbf{x}] = \mu(\mathbf{x}) \approx \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

If this approximation is reasonable, we write the random variable Y conditioned on \mathbf{x} as

$$(Y|\mathbf{X} = \mathbf{x}) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \epsilon$$

ϵ is a random error term.

The Linear Model

Suppose

$$(Y|\mathbf{X} = \mathbf{x}) = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

Consider

$$\mathbf{x}^* = \mathbf{x} + (0, \dots, 0, \underbrace{1}_{j\text{-th entry}}, 0, \dots, 0) = (x_1, x_2, \dots, x_j + 1, \dots, x_k)$$

That is, we change only the j -th entry by 1.

The Linear Model

Suppose

$$(Y|\mathbf{X} = \mathbf{x}) = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

Consider

$$\mathbf{x}^* = \mathbf{x} + (0, \dots, 0, \underbrace{1}_{j\text{-th entry}}, 0, \dots, 0) = (x_1, x_2, \dots, x_j + 1, \dots, x_k)$$

That is, we change only the j -th entry by 1. Then

$$\begin{aligned}(Y|\mathbf{X} = \mathbf{x}^*) &= \beta_0 + \beta_1 x_1 + \cdots + \beta_j (x_j + 1) + \cdots + \beta_k x_k + \epsilon \\ &= \beta_j + \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_j x_j + \cdots + \beta_k x_k}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon\end{aligned}$$

The Linear Model

$$(Y|\mathbf{X} = \mathbf{x}^*) = \underbrace{\beta_j + \mathbb{E}(Y|\mathbf{X} = \mathbf{x})}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x}^*)} + \epsilon$$

Interpretation

The coefficient β_j represents the approximate change in the expected value of Y for a one-unit change in x_j , holding all the other variables constant.

The Linear Model

$$(Y|\mathbf{X} = \mathbf{x}^*) = \underbrace{\beta_j + \mathbb{E}(Y|\mathbf{X} = \mathbf{x})}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x}^*)} + \epsilon$$

Interpretation

The coefficient β_j represents the approximate change in the expected value of Y for a one-unit change in x_j , holding all the other variables constant.

Effect of the feature

This β_j is called the effect of the j -th feature of the expected value of Y .

This β_j effect is the same for all configurations \mathbf{x} .

Is this linear model reasonable?

In the linear model

$$(Y|\mathbf{X} = \mathbf{x}) = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

we simply multiply the j -th feature by a weight β_j and add them up.

It is a simple linear combination of the features.

Isn't it too simple? Yes, it is, but think:

Most times, the effect of a feature may not be linear but it is monotonic:

Is this linear model reasonable?

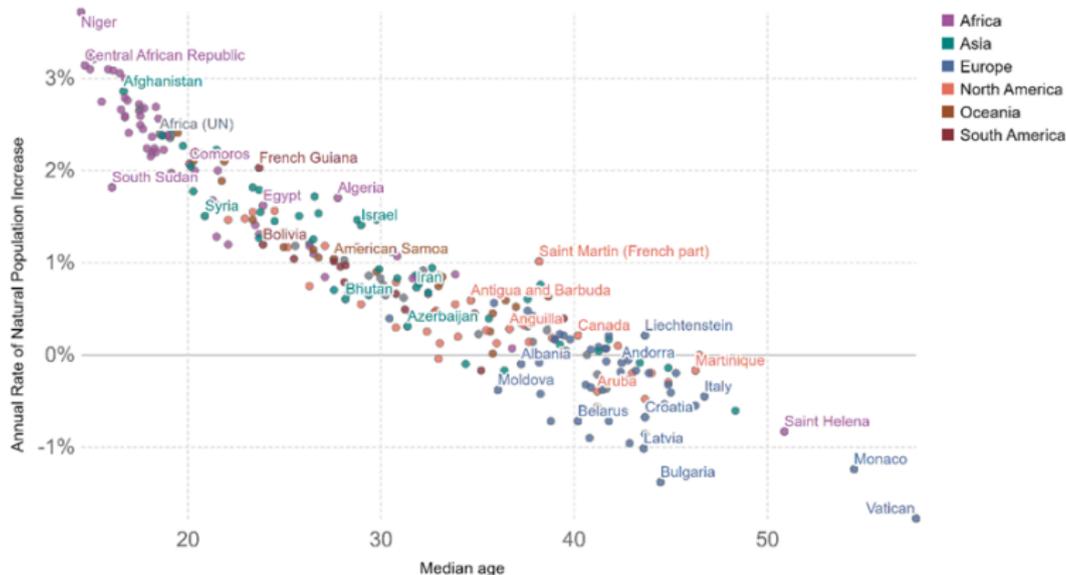
- As we grow old, most health conditions deteriorate (ON AVERAGE) such as mobility, strength, osteoarthritis, etc.
- As the dosage of a toxic substance (e.g., arsenic) increases, the risk of adverse health effects consistently increases.
- Diagnostic indicators (e.g., blood pressure, blood glucose, or Prostate Health Index) deviate further from the normal range, the severity of the associated disease increases.
- The higher the number of cigarettes smoked per day, the higher the risk of coronary heart disease or lung cancer.
- As years of education increase, income generally tends to rise.
- Generally, as the price of a good increases, the quantity demanded decreases.
- The more hours spent studying a subject, the higher the proficiency level achieved.

Monotonic and approximately linear

Population growth rate vs. median age, 2021

Our World
in Data

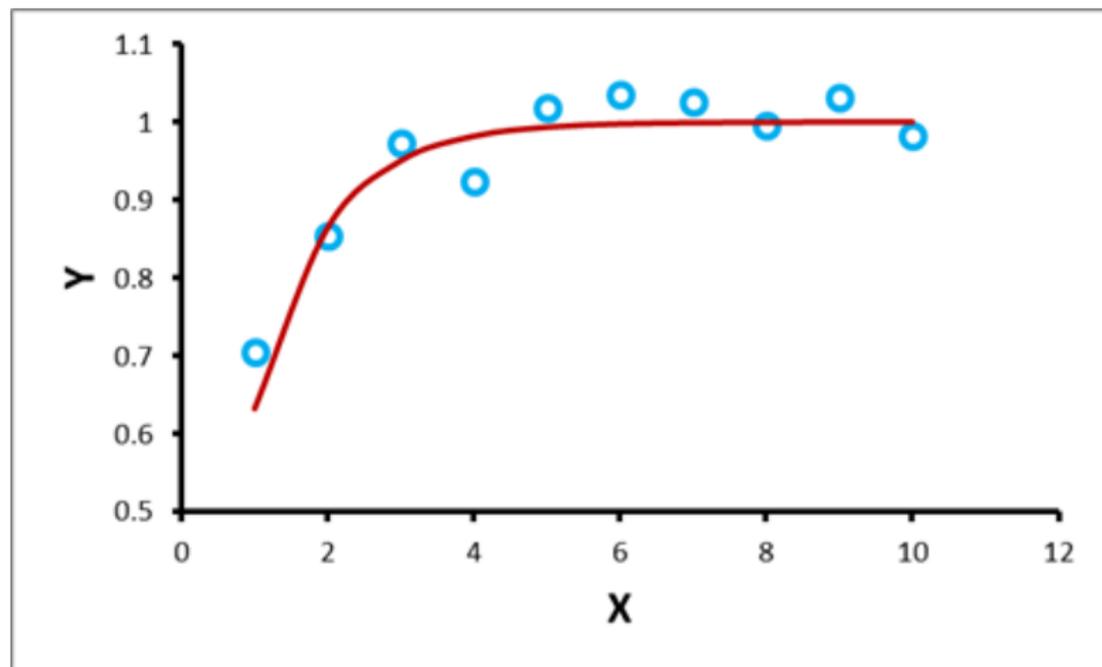
Median age is the age that divides the population into two parts of equal size, that is, there are as many persons with ages above the median as there are with ages below the median. In this metric of population growth, changes due to migration are excluded and only births and deaths are considered.



Source: United Nations, World Population Prospects (2022)

OurWorldInData.org/age-structure • CC BY

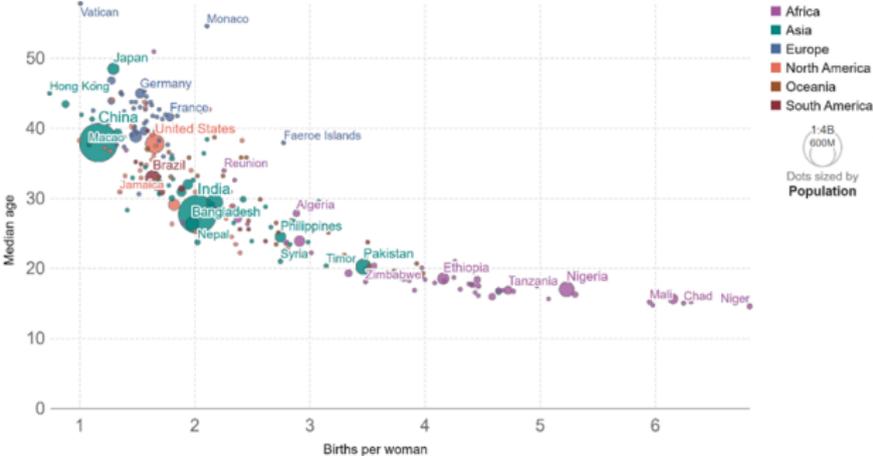
Monotonic but non-linear



Monotonic but non-linear

Median age vs. births per woman, 2021

The median age divides the population in two parts of equal size: that is, there are as many persons with ages above the median as there are with ages below the median.



Source: United Nations, World Population Prospects (2022)

OurWorldInData.org/age-structure • CC BY

Note: The total fertility rate is the number of children that would be born to a woman if she were to live to the end of her child-bearing years and give birth to children at the current age-specific fertility rates.

What if the linear approximation isn't good enough?

Higher-Order Terms

Using a second-order Taylor expansion would introduce:

- Quadratic terms: $\beta_{jj}x_j^2$
- Interaction terms: $\beta_{ij}x_i x_j$ ($i \neq j$)

Model becomes: $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}] \approx \beta_0 + \sum_j \beta_j x_j + \sum_j \beta_{jj} x_j^2 + \sum_{i < j} \beta_{ij} x_i x_j$.

The Linear Regression Model

Summary

The standard linear regression model can be viewed as arising from:

- 1 The goal of predicting Y using \mathbf{X} .
- 2 Recognizing $\mathbb{E}[Y|\mathbf{X}]$ as the optimal MSE predictor.
- 3 Approximating the (potentially complex) function $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ with a first-order Taylor expansion (a linear function).

The Linear Regression Model

Summary

The standard linear regression model can be viewed as arising from:

- 1 The goal of predicting Y using \mathbf{X} .
- 2 Recognizing $\mathbb{E}[Y|\mathbf{X}]$ as the optimal MSE predictor.
- 3 Approximating the (potentially complex) function $\mathbb{E}[Y|\mathbf{X} = \mathbf{x}]$ with a first-order Taylor expansion (a linear function).

Model:

$$Y = \mathbb{E}[Y|\mathbf{X}] + \epsilon \approx (\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k) + \epsilon$$

Or more commonly written as:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \epsilon'$$

where ϵ' includes both the original error ϵ and the approximation error from substituting the true (and unknown) $\mathbb{E}[Y|\mathbf{X}]$ by the linear approximation $\beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k$.

How to estimate the linear model?

We assume a theoretical (approximate) model:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \epsilon$$

How to estimate the linear model?

We assume a theoretical (approximate) model:

$$Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k + \epsilon$$

We have a sample of observed data.

Say, 1500 apartment house prices Y and 30 of their associated characteristics selected because we expect that they have influence on the houses prices: $X_1 =$ total area, $X_2 =$ number of rooms, $X_3 =$ number of suites, $X_4 =$ does it have a gym?, etc.

How to use this data sample to learn (or estimate) the values of the coefficients $\beta_0, \beta_1, \dots, \beta_{30}$ in the linear model

$$Y \text{ (price)} = \beta_0 + \beta_1 X_1 \text{ (area)} + \cdots + \beta_{30} X_{30} + \epsilon ?$$

Organizing the data into matrices

- Prices: a column-vector \mathbf{Y} of dimension 1500.
- The associated features or characteristics: a matrix 1500×30
 - Each row = an apartment house
 - First column (X_1) = Total area
 - Second column (X_2) = house age
 - Third column (X_3) = number of rooms
 - Fourth column (X_4) = number of suites
 - Fifth column (X_5) = number of apartment units on each floor
 - Sixth column (X_6) = number of parking spaces
 - Seventh column (X_7) = Building amenities include a swimming pool? (0 or 1)
 - Etc.
 - Thirty column (X_{30}) = Amenities include gym? (0 or 1)

Vis o matricial

- Prices: a column-vector \mathbf{Y} of dimension 1500. 
- 30 features of the 1500 houses (a matrix of dimension 1500×30) 

$$\mathbf{Y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{1499} \\ y_{1500} \end{pmatrix}$$

$$\begin{pmatrix} \text{area}_1 & \text{age}_1 & \text{rooms}_1 & \cdots & \text{gym}_1 \\ \text{area}_2 & \text{age}_2 & \text{rooms}_2 & \cdots & \text{gym}_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{area}_{1499} & \text{age}_{1499} & \text{rooms}_{1499} & \cdots & \text{gym}_{1499} \\ \text{area}_{1500} & \text{age}_{1500} & \text{rooms}_{1500} & \cdots & \text{gym}_{1500} \end{pmatrix}$$

Linear model implies a weighted sum

- The linear model

$$Y = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_{30}}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

implies that the random price will be predicted by a weighted sum of the 30 features.

Linear model implies a weighted sum

- The linear model

$$Y = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_{30}}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

implies that the random price will be predicted by a weighted sum of the 30 features.

We multiply each of the 30 features by a single constant (the coefficient or weight β_j), sum them up, and add a global constant β_0 .

Linear model implies a weighted sum

- The linear model

$$Y = \underbrace{\beta_0 + \beta_1 x_1 + \cdots + \beta_k x_{30}}_{\mathbb{E}(Y|\mathbf{X}=\mathbf{x})} + \epsilon$$

implies that the random price will be predicted by a weighted sum of the 30 features.

We multiply each of the 30 features by a single constant (the coefficient or weight β_j), sum them up, and add a global constant β_0 .

How do we use the data to learn (or estimate) these coefficients? This is the subject of the next set of slides.