CS70: Lecture 35.

Regression (contd.): Linear and Beyond

CS70: Lecture 35.

Regression (contd.): Linear and Beyond

- 1. Review: Linear Regression (LR), LLSE
- 2. LR: Examples
- 3. Beyond LR: Quadratic Regression
- 4. Conditional Expectation (CE) and properties
- 5. Non-linear Regression: CE = Minimum Mean-Squared Error (MMSE)

Review: Linear Regression – Motivation

Example: 100 people.

Let (X_n, Y_n) = (height, weight) of person *n*, for n = 1, ..., 100:



The blue line is Y = -114.3 + 106.5X. (X in meters, Y in kg.) Best linear fit: Linear Regression.

Review: Covariance Definition

The covariance of X and Y is

$$cov(X,Y) := E[(X - E[X])(Y - E[Y])].$$

Fact

$$cov(X, Y) = E[XY] - E[X]E[Y].$$

Review: Examples of Covariance



Note that E[X] = 0 and E[Y] = 0 in these examples. Then cov(X, Y) = E[XY]. When cov(X, Y) > 0, the RVs X and Y tend to be large or small together. X and Y are said to be positively correlated. When cov(X, Y) < 0, when X is larger, Y tends to be smaller. X and Y are said to be negatively correlated. When cov(X, Y) = 0, we say that X and Y are uncorrelated.

Review: Linear Regression - Non-Bayesian

Definition

Given the samples $\{(X_n, Y_n), n = 1, ..., N\}$, the Linear Regression of *Y* over *X* is

$$\hat{Y} = a + bX$$

where (a, b) minimize

$$\sum_{n=1}^{N}(Y_n-a-bX_n)^2.$$

Thus, $\hat{Y}_n = a + bX_n$ is our guess about Y_n given X_n . The squared error is $(Y_n - \hat{Y}_n)^2$. The LR minimizes the sum of the squared errors. Note: This is a non-Bayesian formulation: there is no prior.

Review: Linear Least Squares Estimate (LLSE)

Definition

Given two RVs X and Y with known distribution Pr[X = x, Y = y], the Linear Least Squares Estimate of Y given X is

$$\hat{Y} = a + bX =: L[Y|X]$$

where (a, b) minimize

$$g(a,b):=E[(Y-a-bX)^2].$$

Thus, $\hat{Y} = a + bX$ is our guess about *Y* given *X*. The squared error is $(Y - \hat{Y})^2$. The LLSE minimizes the expected value of the squared error. Note: This is a Bayesian formulation: there is a prior.

Review: LR: Non-Bayesian or Uniform?

Observe that

$$\frac{1}{N}\sum_{n=1}^{N}(Y_n - a - bX_n)^2 = E[(Y - a - bX)^2]$$

where one assumes that

$$(X, Y) = (X_n, Y_n), \text{ w.p. } \frac{1}{N} \text{ for } n = 1, \dots, N.$$

That is, the non-Bayesian LR is equivalent to the Bayesian LLSE that assumes that (X, Y) is uniform on the set of observed samples.

Thus, we can study the two cases LR and LLSE in one shot. However, the interpretations are different!

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

E[X] =

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_n; \qquad E[Y] =$$

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_n; \qquad E[Y] = \frac{1}{N} \sum_{n=1}^{N} Y_n$$

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_n; \qquad E[Y] = \frac{1}{N} \sum_{n=1}^{N} Y_n$$

 $Var[X] = E[X^2] - (E[X])^2 =$

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_n; \qquad E[Y] = \frac{1}{N} \sum_{n=1}^{N} Y_n$$
$$Var[X] = E[X^2] - (E[X])^2 = \frac{1}{N} \sum_{n=1}^{N} (X_n)^2 - (\frac{1}{N} \sum_{n=1}^{N} (X_n))^2$$

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_n; \qquad E[Y] = \frac{1}{N} \sum_{n=1}^{N} Y_n$$
$$Var[X] = E[X^2] - (E[X])^2 = \frac{1}{N} \sum_{n=1}^{N} (X_n)^2 - (\frac{1}{N} \sum_{n=1}^{N} (X_n))^2$$

Cov(X, Y) = E[XY] - E[X]E[Y]

Theorem

Consider two RVs X, Y with a given distribution Pr[X = x, Y = y]. Then, $L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X, Y)}{var(X)}(X - E[X]).$

Non-Bayesian setting:

$$E[X] = \frac{1}{N} \sum_{n=1}^{N} X_{n}; \qquad E[Y] = \frac{1}{N} \sum_{n=1}^{N} Y_{n}$$
$$Var[X] = E[X^{2}] - (E[X])^{2} = \frac{1}{N} \sum_{n=1}^{N} (X_{n})^{2} - (\frac{1}{N} \sum_{n=1}^{N} (X_{n}))^{2}$$
$$Cov(X, Y) = E[XY] - E[X]E[Y] = \frac{1}{N} \sum_{n=1}^{N} (X_{n}Y_{n}) - (\frac{1}{N} \sum_{n=1}^{N} X_{n})(\frac{1}{N} \sum_{n=1}^{N} Y_{n})$$

LR: Illustration



LR: Illustration



Note that

▶ the LR line goes through (*E*[*X*], *E*[*Y*])

LR: Illustration



Note that

▶ the LR line goes through (*E*[*X*], *E*[*Y*])

• its slope is
$$\frac{cov(X,Y)}{var(X)}$$
.



Example : "Removing noise or de-noising" Y: temp. in a room (quantity of interest) $\left[\Upsilon = \mathcal{N}\left(\mathcal{M}_{r}, \sigma_{Y}^{2}\right)\right]$ Z: thermal noise of temp. sensor $\left[Z = \mathcal{N}(0, \sigma_{z}^{2}) \right]$ X: observed (noisy) temp measurement @ server Z(noise) (true temp.) X = Y + Z(observed.) L[Y|X] = a+bx + E (est. +own) L[Y|X]=Y=a+bX, where a,b chosen to min $IE\left[\frac{\text{rest.}}{\text{error}}\right]^2 = IE\left[(Y-Y)^2\right]$ $\underline{LLSE}: \quad \widehat{\mathbf{Y}} = \underline{IE}[\mathbf{Y}] + \frac{\operatorname{cov}(\mathbf{X},\mathbf{Y})}{\operatorname{var}(\mathbf{X})} [\mathbf{X} - IE[\mathbf{X}]]$

 $\cdot IE[X] = IE[Y + Z] = IE[Y] + E[Z] = \mu_{Y}$ $\cdot cov(X,Y) = IE[XY] - IE[X] \cdot IE[Y]$ $= IE[(Y+2)Y] - IE[Y+2] \cdot IE[Y]$ IE(Y2]+E[YZ] σy2+μy2+ [E[Y]·E[] $\Rightarrow cov(X_1Y) = \sigma_Y^2 + M_Y^2 - M_Y^2 = \sigma_Y^2$ $var(X) = var(Y) + var(Z) = (\sigma_Y^{2+\sigma_Z^2})$ $\hat{Y} = L[Y|X] = \mu_Y + \left(\frac{\sigma_Y^2}{\sigma_Y^2 + \sigma_{\Xi^2}}\right)(X - \mu_Y)$ Remarks: U) If of 20 (no nouse) ⇒ Y = X ("believe) (2) $If G_{\overline{z}}^{2} > 5G_{Y}^{2}(v, noisy) \Rightarrow \widehat{Y} \cong \mu_{Y} = IE[Y]$ ("believe the, madel & not madel & not (a) If $\mu_{Y}=70^{\circ}F$, $\sigma_{Y}^{2}=5$, $\sigma_{\Xi}^{2}=2$ $\hat{Y}=70+\underline{5}(X-70)$





We find:

E[X] =



We find:

E[X] = 0;



We find:

E[X] = 0; E[Y] =



We find:

E[X] = 0; E[Y] = 0;



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] =$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2;$$



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] =$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = 1/2;$$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = 1/2;$$

 $var[X] = E[X^2] - E[X]^2 =$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = 1/2;$$

 $var[X] = E[X^2] - E[X]^2 = 1/2;$



$$E[X] = 0; E[Y] = 0; E[X2] = 1/2; E[XY] = 1/2;$$

var[X] = E[X²] - E[X]² = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] =



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = 1/2;$$

 $var[X] = E[X^2] - E[X]^2 = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = 1/2;$



$$E[X] = 0; E[Y] = 0; E[X^{2}] = 1/2; E[XY] = 1/2;$$

$$var[X] = E[X^{2}] - E[X]^{2} = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = 1/2;$$

$$LR: \hat{Y} = E[Y] + \frac{cov(X, Y)}{var[X]}(X - E[X]) =$$


$$E[X] = 0; E[Y] = 0; E[X^{2}] = 1/2; E[XY] = 1/2;$$

$$var[X] = E[X^{2}] - E[X]^{2} = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = 1/2;$$

$$LR: \hat{Y} = E[Y] + \frac{cov(X, Y)}{var[X]}(X - E[X]) = X.$$





We find:

E[X] =



We find:

E[X] = 0;



We find:

E[X] = 0; E[Y] =



We find:

E[X] = 0; E[Y] = 0;



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] =$



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] = 1/2;$



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] =$



We find:

 $E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = -1/2;$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = -1/2;$$

 $var[X] = E[X^2] - E[X]^2 =$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = -1/2;$$

 $var[X] = E[X^2] - E[X]^2 = 1/2;$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = -1/2;$$

 $var[X] = E[X^2] - E[X]^2 = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] =$



$$E[X] = 0; E[Y] = 0; E[X^2] = 1/2; E[XY] = -1/2;$$

 $var[X] = E[X^2] - E[X]^2 = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = -1/2;$



$$E[X] = 0; E[Y] = 0; E[X^{2}] = 1/2; E[XY] = -1/2;$$

$$var[X] = E[X^{2}] - E[X]^{2} = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = -1/2;$$

$$LR: \hat{Y} = E[Y] + \frac{cov(X, Y)}{var[X]}(X - E[X]) =$$



$$E[X] = 0; E[Y] = 0; E[X^{2}] = 1/2; E[XY] = -1/2;$$

$$var[X] = E[X^{2}] - E[X]^{2} = 1/2; cov(X, Y) = E[XY] - E[X]E[Y] = -1/2;$$

$$LR: \hat{Y} = E[Y] + \frac{cov(X, Y)}{var[X]}(X - E[X]) = -X.$$

We saw that the LLSE of Y given X is

$$L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X]).$$

We saw that the LLSE of Y given X is

$$L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X]).$$

How good is this estimator?

We saw that the LLSE of Y given X is

$$L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X]).$$

How good is this estimator? That is, what is the mean squared estimation error?

We saw that the LLSE of Y given X is

$$L[Y|X] = \hat{Y} = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X]).$$

How good is this estimator? That is, what is the mean squared estimation error?

We find

$$\begin{split} & E[|Y - L[Y|X]|^2] = E[(Y - E[Y] - (cov(X, Y)/var(X))(X - E[X]))^2] \\ &= E[(Y - E[Y])^2] - 2(cov(X, Y)/var(X))E[(Y - E[Y])(X - E[X])] \\ &+ (cov(X, Y)/var(X))^2E[(X - E[X])^2 \\ &= var(Y) - \frac{cov(X, Y)^2}{var(X)}. \end{split}$$

Without observations, the estimate is E[Y] = 0. The error is var(Y). Observing X reduces the error.

Linear Regression

Linear Regression

1. Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X])$

Linear Regression

- 1. Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X E[X])$
- 2. Non-Bayesian: minimize $\sum_{n} (Y_n a bX_n)^2$

Linear Regression

- 1. Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X E[X])$
- 2. Non-Bayesian: minimize $\sum_{n} (Y_n a bX_n)^2$
- 3. Bayesian: minimize $E[(Y-a-bX)^2]$

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is?

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

LR: Restriction to linear functions: g(X) = a + bX.

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

LR: Restriction to linear functions: g(X) = a + bX.

With no such constraints, what is the best g(X)?

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

LR: Restriction to linear functions: g(X) = a + bX.

With no such constraints, what is the best g(X)?

Answer:

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

LR: Restriction to linear functions: g(X) = a + bX.

With no such constraints, what is the best g(X)?

Answer: E[Y|X].

Goal: guess the value of Y in the expected squared error sense. We know nothing about Y other than its distribution. Our best guess is? E[Y].

Now assume we make some observation X related to Y.

How do we use that observation to improve our guess about Y?

Idea: use a function g(X) of the observation to estimate *Y*.

LR: Restriction to linear functions: g(X) = a + bX.

With no such constraints, what is the best g(X)?

Answer: E[Y|X].

This is called the Conditional Expectation (CE).

Nonlinear Regression: Motivation
There are many situations where a good guess about Y given X is not linear.

There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight),

There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight), (school years, income),

There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight), (school years, income), (PSA level, cancer risk).

There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight), (school years, income), (PSA level, cancer risk).



There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight), (school years, income), (PSA level, cancer risk).



Our goal:

There are many situations where a good guess about Y given X is not linear.

E.g., (diameter of object, weight), (school years, income), (PSA level, cancer risk).



Our goal: explore estimates $\hat{Y} = g(X)$ for nonlinear functions $g(\cdot)$.

Let X, Y be two random variables defined on the same probability space.

Let X, Y be two random variables defined on the same probability space.

Definition:

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$. Derivation:

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*.

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*. We get

$$0 = E[Y-a-bX-cX^2]$$

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of *Y* over *X* is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*. We get

$$0 = E[Y-a-bX-cX^{2}]$$

$$0 = E[(Y-a-bX-cX^{2})X$$

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of Y over X is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*. We get

$$0 = E[Y - a - bX - cX^{2}]$$

$$0 = E[(Y - a - bX - cX^{2})X]$$

$$0 = E[(Y - a - bX - cX^{2})X^{2}]$$

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of *Y* over *X* is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*. We get

$$0 = E[Y - a - bX - cX^{2}]$$

$$0 = E[(Y - a - bX - cX^{2})X]$$

$$0 = E[(Y - a - bX - cX^{2})X^{2}]$$

We solve these three equations in the three unknowns (a, b, c).

Let X, Y be two random variables defined on the same probability space.

Definition: The quadratic regression of *Y* over *X* is the random variable

$$Q[Y|X] = a + bX + cX^2$$

where a, b, c are chosen to minimize $E[(Y - a - bX - cX^2)^2]$.

Derivation: We set to zero the derivatives w.r.t. *a*, *b*, *c*. We get

$$0 = E[Y - a - bX - cX^{2}]$$

$$0 = E[(Y - a - bX - cX^{2})X]$$

$$0 = E[(Y - a - bX - cX^{2})X^{2}]$$

We solve these three equations in the three unknowns (a, b, c).

Conditional Expectation

Definition Let *X* and *Y* be RVs on Ω .

Definition Let *X* and *Y* be RVs on Ω . The conditional expectation of *Y* given *X* is defined as

E[Y|X] = g(X)

Definition Let *X* and *Y* be RVs on Ω . The conditional expectation of *Y* given *X* is defined as

E[Y|X] = g(X)

where

$$g(x) := E[Y|X = x] := \sum_{y} yPr[Y = y|X = x].$$

Definition Let *X* and *Y* be RVs on Ω . The conditional expectation of *Y* given *X* is defined as

E[Y|X] = g(X)

where

$$g(x) := E[Y|X = x] := \sum_{y} yPr[Y = y|X = x].$$

Have we seen this before?

Have we seen this before? Yes.

Have we seen this before? Yes.

Is anything new?

Have we seen this before? Yes.

Is anything new? Yes.

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Have we seen this before? Yes.

```
Is anything new? Yes.
```

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal?

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite!

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of *X*.

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of X.

This is similar: E[Y|X] = g(X) for some function $g(\cdot)$.

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of X.

This is similar: E[Y|X] = g(X) for some function $g(\cdot)$.

In general, g(X) is not linear, i.e., not a+bX.

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of X.

This is similar: E[Y|X] = g(X) for some function $g(\cdot)$.

In general, g(X) is not linear, i.e., not a+bX. It could be that $g(X) = a+bX+cX^2$.
Deja vu, all over again?

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of X.

This is similar: E[Y|X] = g(X) for some function $g(\cdot)$.

In general, g(X) is not linear, i.e., not a + bX. It could be that $g(X) = a + bX + cX^2$. Or that $g(X) = 2\sin(4X) + \exp\{-3X\}$.

Deja vu, all over again?

Have we seen this before? Yes.

Is anything new? Yes.

The idea of defining g(x) = E[Y|X = x] and then E[Y|X] = g(X).

Big deal? Quite! Simple but most convenient.

Recall that L[Y|X] = a + bX is a function of X.

This is similar: E[Y|X] = g(X) for some function $g(\cdot)$.

In general, g(X) is not linear, i.e., not a + bX. It could be that $g(X) = a + bX + cX^2$. Or that $g(X) = 2\sin(4X) + \exp\{-3X\}$. Or something else.

$$E[Y|X=x] = \sum_{y} y Pr[Y=y|X=x]$$

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

(a) X, Y independent $\Rightarrow E[Y|X] = E[Y];$

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

- (a) X, Y independent $\Rightarrow E[Y|X] = E[Y];$
- (b) E[aY+bZ|X] = aE[Y|X]+bE[Z|X];

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

(a) X, Y independent $\Rightarrow E[Y|X] = E[Y];$

- (b) E[aY+bZ|X] = aE[Y|X]+bE[Z|X];
- (c) $E[Yh(X)|X] = h(X)E[Y|X], \forall h(\cdot);$

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

(a) X, Y independent $\Rightarrow E[Y|X] = E[Y];$

- (b) E[aY+bZ|X] = aE[Y|X] + bE[Z|X];
- (c) $E[Yh(X)|X] = h(X)E[Y|X], \forall h(\cdot);$

(d) E[E[Y|X]] = E[Y].

$$E[Y|X=x] = \sum_{y} yPr[Y=y|X=x]$$

Theorem

(a) X, Y independent $\Rightarrow E[Y|X] = E[Y];$

- (b) E[aY+bZ|X] = aE[Y|X] + bE[Z|X];
- (c) $E[Yh(X)|X] = h(X)E[Y|X], \forall h(\cdot);$

(d) E[E[Y|X]] = E[Y].

Let X, Y, Z be i.i.d. with mean 0 and variance 1.

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

$$E[2+5X+7XY+11X^2+13X^3Z^2|X]$$

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

$$E[2+5X+7XY+11X^{2}+13X^{3}Z^{2}|X]$$

= 2+5X+7XE[Y|X]+11X^{2}+13X^{3}E[Z^{2}|X]

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

$$E[2+5X+7XY+11X^{2}+13X^{3}Z^{2}|X]$$

= 2+5X+7XE[Y|X]+11X^{2}+13X^{3}E[Z^{2}|X]
= 2+5X+7XE[Y]+11X^{2}+13X^{3}E[Z^{2}]

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

$$E[2+5X+7XY+11X^{2}+13X^{3}Z^{2}|X]$$

= 2+5X+7XE[Y|X]+11X^{2}+13X^{3}E[Z^{2}|X]
= 2+5X+7XE[Y]+11X^{2}+13X^{3}E[Z^{2}]
= 2+5X+11X^{2}+13X^{3}(var[Z]+E[Z]^{2})

Let X, Y, Z be i.i.d. with mean 0 and variance 1. We want to calculate

$$E[2+5X+7XY+11X^2+13X^3Z^2|X].$$

$$\begin{split} & E[2+5X+7XY+11X^2+13X^3Z^2|X] \\ &= 2+5X+7XE[Y|X]+11X^2+13X^3E[Z^2|X] \\ &= 2+5X+7XE[Y]+11X^2+13X^3E[Z^2] \\ &= 2+5X+11X^2+13X^3(var[Z]+E[Z]^2) \\ &= 2+5X+11X^2+13X^3. \end{split}$$

(Conditional Expectation = Minimum Mean Squared Error) **Theorem**

g(X) := E[Y|X] is the function of X that minimizes $E[(Y - g(X))^2]$.

(Conditional Expectation = Minimum Mean Squared Error) **Theorem**

g(X) := E[Y|X] is the function of X that minimizes $E[(Y - g(X))^2]$.

That is, E[Y|X] is the 'best' guess about Y based on X.

(Conditional Expectation = Minimum Mean Squared Error) **Theorem**

g(X) := E[Y|X] is the function of X that minimizes $E[(Y - g(X))^2]$.

That is, E[Y|X] is the 'best' guess about Y based on X.

Specifically, it is the function g(X) of X that

minimizes $E[(Y - g(X))^2]$.

(Conditional Expectation = Minimum Mean Squared Error) **Theorem**

g(X) := E[Y|X] is the function of X that minimizes $E[(Y - g(X))^2]$.

That is, E[Y|X] is the 'best' guess about Y based on X.

Specifically, it is the function g(X) of X that

minimizes $E[(Y - g(X))^2]$.





Linear and Non-Linear Regression: Conditional Expectation

Summary

Linear and Non-Linear Regression: Conditional Expectation

► Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X - E[X])$

Summary

Linear and Non-Linear Regression: Conditional Expectation

- ► Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X E[X])$
- ► Non-linear Regression: MMSE: E[Y|X] minimizes E[(Y - g(X))²] over all g(·)
- Definition: $E[Y|X] := \sum_{y} y Pr[Y = y|X = x]$

Summary

Linear and Non-Linear Regression: Conditional Expectation

- ► Linear Regression: $L[Y|X] = E[Y] + \frac{cov(X,Y)}{var(X)}(X E[X])$
- ► Non-linear Regression: MMSE: E[Y|X] minimizes E[(Y - g(X))²] over all g(·)
- Definition: $E[Y|X] := \sum_{y} y Pr[Y = y|X = x]$