# YSU ASDS, Statistics, Fall 2019
## Lecture 28

Michael Poghosyan

02 Dec 2019

# Contents

# Last Lecture ReCap

- Give two alternative definitions of the $p$-Value of a Test.

# Last Lecture ReCap

- Give two alternative definitions of the $p$-Value of a Test.
- Using $p$-Values, in which case we Reject Null?

# Linear Regression

# Linear Regression

From the Statistical Learning Perspective

# Intro to Linear Regression

**Note:** From this slide on, **LR** will mean Linear Regression, not the Likelihood Ratio ☺

# Intro to Linear Regression

**Note:** From this slide on, **LR** will mean Linear Regression, not the Likelihood Ratio ‿

So far, except some exceptions ‿, we have considered inference about one r.v. and the Distribution behind: our r.v. was $X \sim \mathcal{F}_\theta$, and our aim was, given a Random Sample

$$X_1, X_2, ..., X_n \sim \mathcal{F}_\theta,$$

to get some information about $\theta$, in particular,

# Intro to Linear Regression

**Note:** From this slide on, **LR** will mean Linear Regression, not the Likelihood Ratio ⌣

So far, except some exceptions ⌣, we have considered inference about one r.v. and the Distribution behind: our r.v. was $X \sim \mathcal{F}_\theta$, and our aim was, given a Random Sample

$$X_1, X_2, ..., X_n \sim \mathcal{F}_\theta,$$

to get some information about $\theta$, in particular,

▶ to find a good Point Estimator and Estimate;

# Intro to Linear Regression

**Note:** From this slide on, **LR** will mean Linear Regression, not the Likelihood Ratio ‿

So far, except some exceptions ‿, we have considered inference about one r.v. and the Distribution behind: our r.v. was $X \sim \mathcal{F}_\theta$, and our aim was, given a Random Sample

$$X_1, X_2, ..., X_n \sim \mathcal{F}_\theta,$$

to get some information about $\theta$, in particular,

- ▶ to find a good Point Estimator and Estimate;
- ▶ to find a CI for $\theta$ of given CL;

# Intro to Linear Regression

**Note:** From this slide on, **LR** will mean Linear Regression, not the Likelihood Ratio ⌣

So far, except some exceptions ⌣, we have considered inference about one r.v. and the Distribution behind: our r.v. was $X \sim \mathcal{F}_\theta$, and our aim was, given a Random Sample

$$X_1, X_2, ..., X_n \sim \mathcal{F}_\theta,$$

to get some information about $\theta$, in particular,

- ▶ to find a good Point Estimator and Estimate;

- ▶ to find a CI for $\theta$ of given CL;

- ▶ to Test a Hypothesis about $\theta$, say, is it likely that $\theta = 3.1415$ or not.

# Intro to Linear Regression

Now we want to talk about Modeling and inference for two (or more) Dependent Random variables, to talk about Modeling the Relationship between two or more Random Variables.

# Intro to Linear Regression

Now we want to talk about Modeling and inference for two (or more) Dependent Random variables, to talk about Modeling the Relationship between two or more Random Variables.

Recall that, in the Descriptive Statistics part, we considered two Datasets, and defined the Sample Covariance and Correlation Coefficients, to measure the Linear Relationship between that Datasets.

# Intro to Linear Regression

Now we want to talk about Modeling and inference for two (or more) Dependent Random variables, to talk about Modeling the Relationship between two or more Random Variables.

Recall that, in the Descriptive Statistics part, we considered two Datasets, and defined the Sample Covariance and Correlation Coefficients, to measure the Linear Relationship between that Datasets. That was defined for two **Numerical Dataset**, without any assumptions behind the Process generating that Datasets.

# Intro to Linear Regression

Now we want to talk about Modeling and inference for two (or more) Dependent Random variables, to talk about Modeling the Relationship between two or more Random Variables.

Recall that, in the Descriptive Statistics part, we considered two Datasets, and defined the Sample Covariance and Correlation Coefficients, to measure the Linear Relationship between that Datasets. That was defined for two **Numerical Dataset**, without any assumptions behind the Process generating that Datasets. Now, if we assume that that Datasets are coming from some Distribution, we are at the stage of doing a Statistical Inference, Statistical Analysis.

# Elements of the Supervised Learning

# Supervised Learning

So we start again from the Data.

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

► The Input Space $\mathcal{X}$

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

- The Input Space $\mathcal{X}$
- The Output Space $\mathcal{Y}$

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

- The Input Space $\mathcal{X}$
- The Output Space $\mathcal{Y}$

We will assume $\mathcal{X} \subset \mathbb{R}^d$ (or, maybe, in other $d$-Dim Space), and a typical element $\mathbf{x}$ of $\mathcal{X}$ will have the form

$$\mathbf{x} = (x_1, x_2, ..., x_d)$$

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

- ▶ The Input Space $\mathcal{X}$
- ▶ The Output Space $\mathcal{Y}$

We will assume $\mathcal{X} \subset \mathbb{R}^d$ (or, maybe, in other $d$-Dim Space), and a typical element $\mathbf{x}$ of $\mathcal{X}$ will have the form

$$\mathbf{x} = (x_1, x_2, ..., x_d)$$

We will call $x_k$-s to be the **Features** of $\mathbf{x}$.

# Supervised Learning

So we start again from the Data.

Here we assume we are given:

- ▶ The Input Space $\mathcal{X}$
- ▶ The Output Space $\mathcal{Y}$

We will assume $\mathcal{X} \subset \mathbb{R}^d$ (or, maybe, in other $d$-Dim Space), and a typical element $\mathbf{x}$ of $\mathcal{X}$ will have the form

$$\mathbf{x} = (x_1, x_2, ..., x_d)$$

We will call $x_k$-s to be the **Features** of $\mathbf{x}$.

We will assume also that $\mathcal{Y} \subset \mathbb{R}$, and we will call the elements of $\mathcal{Y}$ to be the **Labels**.

# Supervised Learning

In the Supervised Learning Problems, we have a Data of the form:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n).$$

# Supervised Learning

In the Supervised Learning Problems, we have a Data of the form:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n).$$

Here we interpret $\mathbf{x}_k$ and $y_k$ as the Feature vector and the Label of the Observation (Object) $k$.

# Supervised Learning

In the Supervised Learning Problems, we have a Data of the form:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n).$$

Here we interpret $\mathbf{x}_k$ and $y_k$ as the Feature vector and the Label of the Observation (Object) $k$.

So we know the labels of our $n$ Observations.

# Supervised Learning

In the Supervised Learning Problems, we have a Data of the form:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ..., (\mathbf{x}_n, y_n).$$

Here we interpret $\mathbf{x}_k$ and $y_k$ as the Feature vector and the Label of the Observation (Object) $k$.

So we know the labels of our $n$ Observations.

**Problem:** Given a Feature vector $\mathbf{x}$, other than $\mathbf{x}_k$, predict its Label $y$.

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

The Features $x_k$ can be:

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

The Features $x_k$ can be:

- Binary, if $x_k \in \{0, 1\}$ or $x_k \in \{-1, 1\}$

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

The Features $x_k$ can be:

- Binary, if $x_k \in \{0, 1\}$ or $x_k \in \{-1, 1\}$
- Nominal/Categorical, if the set of possible values of $x_k$ is finite, and no intrinsic order exists in that set

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

The Features $x_k$ can be:

- ▶ Binary, if $x_k \in \{0, 1\}$ or $x_k \in \{-1, 1\}$
- ▶ Nominal/Categorical, if the set of possible values of $x_k$ is finite, and no intrinsic order exists in that set
- ▶ Ordinal, if the set of possible values of $x_k$ is finite, and there is a natural order in that set

# Features and Labels

Recall that our Feature Vector $\mathbf{x} \in \mathcal{X}$ has the form:

$$\mathbf{x} = (x_1, x_2, ..., x_d).$$

The Features $x_k$ can be:

- ▶ Binary, if $x_k \in \{0, 1\}$ or $x_k \in \{-1, 1\}$
- ▶ Nominal/Categorical, if the set of possible values of $x_k$ is finite, and no intrinsic order exists in that set
- ▶ Ordinal, if the set of possible values of $x_k$ is finite, and there is a natural order in that set
- ▶ Numerical/Quantitative, if $x_k \in \mathbb{R}$

# Features and Labels

The Labels can be:

# Features and Labels

The Labels can be:

**Classification Problems:**

- $\mathcal{Y} = \{-1, 1\}$ or $\mathcal{Y} = \{0, 1\}$ - **Binary Classification**

# Features and Labels

The Labels can be:

**Classification Problems:**

- $\mathcal{Y} = \{-1, 1\}$ or $\mathcal{Y} = \{0, 1\}$ - **Binary Classification**
- $\mathcal{Y} = \{1, 2, ..., K\}$ - $K$-**class Classification**

# Features and Labels

The Labels can be:

**Classification Problems:**

- $\mathcal{Y} = \{-1, 1\}$ or $\mathcal{Y} = \{0, 1\}$ - **Binary Classification**
- $\mathcal{Y} = \{1, 2, ..., K\}$ - *K*-**class Classification**

**Regression Problems:**

- $\mathcal{Y} = \mathbb{R}$ - **1D Regression**

# Supervised Learning

So, having a Dataset of Observations with Labels, we want to predict the Label for a new Observation. Of course, we cannot do this unless we will assume there is some structure, some relationship in the Data.

# Supervised Learning

So, having a Dataset of Observations with Labels, we want to predict the Label for a new Observation. Of course, we cannot do this unless we will assume there is some structure, some relationship in the Data. Mathematically, we will assume that behind our Data we have a Probability Distribution, generating our Data.

# Supervised Learning

So, having a Dataset of Observations with Labels, we want to predict the Label for a new Observation. Of course, we cannot do this unless we will assume there is some structure, some relationship in the Data. Mathematically, we will assume that behind our Data we have a Probability Distribution, generating our Data.

We will assume that the observation $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is a realization of the pair of r.v.s $(\mathbf{X}, Y)$ that is coming from some unknown Distribution $\mathcal{F}$:

$$(\mathbf{X}, Y) \sim \mathcal{F}.$$

# Supervised Learning

So, having a Dataset of Observations with Labels, we want to predict the Label for a new Observation. Of course, we cannot do this unless we will assume there is some structure, some relationship in the Data. Mathematically, we will assume that behind our Data we have a Probability Distribution, generating our Data.

We will assume that the observation $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is a realization of the pair of r.v.s $(\mathbf{X}, Y)$ that is coming from some unknown Distribution $\mathcal{F}$:

$$(\mathbf{X}, Y) \sim \mathcal{F}.$$

So we "encode" our Data $(\mathbf{x}_k, y_k)$ as being a realisation of a r.v. $(\mathbf{X}_k, Y_k)$.

# Supervised Learning

So, having a Dataset of Observations with Labels, we want to predict the Label for a new Observation. Of course, we cannot do this unless we will assume there is some structure, some relationship in the Data. Mathematically, we will assume that behind our Data we have a Probability Distribution, generating our Data.

We will assume that the observation $(\mathbf{x}, y) \in \mathcal{X} \times \mathcal{Y}$ is a realization of the pair of r.v.s $(\mathbf{X}, Y)$ that is coming from some unknown Distribution $\mathcal{F}$:

$$(\mathbf{X}, Y) \sim \mathcal{F}.$$

So we "encode" our Data $(\mathbf{x}_k, y_k)$ as being a realisation of a r.v. $(\mathbf{X}_k, Y_k)$.

The general idea/Problem is, having Data, to infer $\mathcal{F}$.

# Supervised Learning

This is a very general problem, so we consider the following:

# Supervised Learning

This is a very general problem, so we consider the following: Given a sequence of IID r.v.s

$$(\mathbf{X}_1, Y_1), (\mathbf{X}_2, Y_2), ..., (\mathbf{X}_n, Y_n)$$

construct a "good" Prediction Function

$$g : \mathcal{X} \to \mathcal{Y},$$

that will predict the Label of **X**.

# Supervised Learning, Loss Function

Here we need to talk about different things:

# Supervised Learning, Loss Function

Here we need to talk about different things:

▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

# Supervised Learning, Loss Function

Here we need to talk about different things:

▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

▶ How to construct good Predictors?

# Supervised Learning, Loss Function

Here we need to talk about different things:

- ▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

- ▶ How to construct good Predictors?

Construction of $g$, using the Data we have, is called **Training**.

# Supervised Learning, Loss Function

Here we need to talk about different things:

- ▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

- ▶ How to construct good Predictors?

Construction of $g$, using the Data we have, is called **Training**.
Predicting the values for new observations is called **Testing**.

# Supervised Learning, Loss Function

Here we need to talk about different things:

- ▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

- ▶ How to construct good Predictors?

Construction of $g$, using the Data we have, is called **Training**.
Predicting the values for new observations is called **Testing**.

To assess goodness of the Predictor $g$, we take a **Loss** function.

# Supervised Learning, Loss Function

Here we need to talk about different things:

▶ What is the meaning that $g$ is giving "good" labels, is a "good" Predictor, how to assess that?

▶ How to construct good Predictors?

Construction of $g$, using the Data we have, is called **Training**. Predicting the values for new observations is called **Testing**.

To assess goodness of the Predictor $g$, we take a **Loss** function. We will call any function of the form

$$\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$$

a **Loss** function, and we will assume that:

$$\ell(y_1, y_2) \geq 0, \qquad \forall y_1, y_2 \in \mathcal{Y}, \qquad \text{and} \qquad \ell(y, y) = 0.$$

# Loss Function

Some known Loss Functions are:

# Loss Function

Some known Loss Functions are:

▶ For The Binary Classification:

$$\ell(y_1, y_2) = \left\{ \begin{array}{ll} 1, & y_1 \neq y_2 \\ 0, & y_1 = y_2 \end{array} \right.$$

---

[1]Here we use the Indicator, $\mathbf{1}(x)$ function with $\mathbf{1}(\textit{True}) = 1$ and $\mathbf{1}(\textit{False}) = 0$.

# Loss Function

Some known Loss Functions are:

▶ For The Binary Classification:

$$\ell(y_1, y_2) = \begin{cases} 1, & y_1 \neq y_2 \\ 0, & y_1 = y_2 \end{cases}$$

We denote this[1] as $\mathbf{1}(y_1 \neq y_2)$.

---

[1]Here we use the Indicator, $\mathbf{1}(x)$ function with $\mathbf{1}(True) = 1$ and $\mathbf{1}(False) = 0$.

# Loss Function

Some known Loss Functions are:

▶ For The Binary Classification:

$$\ell(y_1, y_2) = \left\{ \begin{array}{ll} 1, & y_1 \neq y_2 \\ 0, & y_1 = y_2 \end{array} \right.$$

We denote this[1] as $\mathbf{1}(y_1 \neq y_2)$.

▶ For 1D Regression:

▶ $\ell(y_1, y_2) = (y_1 - y_2)^2$ - Quadratic Loss;

---

[1]Here we use the Indicator, $\mathbf{1}(x)$ function with $\mathbf{1}(True) = 1$ and $\mathbf{1}(False) = 0$.

# Loss Function

Some known Loss Functions are:

▶ For The Binary Classification:

$$\ell(y_1, y_2) = \begin{cases} 1, & y_1 \neq y_2 \\ 0, & y_1 = y_2 \end{cases}$$

We denote this[1] as $\mathbf{1}(y_1 \neq y_2)$.

▶ For 1D Regression:

  ▶ $\ell(y_1, y_2) = (y_1 - y_2)^2$ - Quadratic Loss;
  ▶ $\ell(y_1, y_2) = |y_1 - y_2|$ - Absolute Error Loss;

---

[1]Here we use the Indicator, $\mathbf{1}(x)$ function with $\mathbf{1}(True) = 1$ and $\mathbf{1}(False) = 0$.

## Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

## Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$.

# Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$.

# Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$. Then the Loss incurring will be

$$\ell(Y, g(\mathbf{X})).$$

## Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$. Then the Loss incurring will be

$$\ell(Y, g(\mathbf{X})).$$

We want to have this Loss as small as possible.

## Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$. Then the Loss incurring will be

$$\ell(Y, g(\mathbf{X})).$$

We want to have this Loss as small as possible. Well, under our setting, this will be a R.V., so we define the **Average Loss** or the **Risk** of the Predictor $g$ to be

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

# Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$. Then the Loss incurring will be

$$\ell(Y, g(\mathbf{X})).$$

We want to have this Loss as small as possible. Well, under our setting, this will be a R.V., so we define the **Average Loss** or the **Risk** of the Predictor $g$ to be

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

Here the Expectation is over the Distribution of $(\mathbf{X}, Y)$, i.e., $\mathcal{F}$.

# Learning Problem

Now assume we have a Predictor $g : \mathcal{X} \to \mathcal{Y}$, and a Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$.

Assume $(\mathbf{X}, Y) \sim \mathcal{F}$. So $\mathbf{X}$ is our Feature Vector, and we Predict its label as $g(\mathbf{X})$. Then the Loss incurring will be

$$\ell(Y, g(\mathbf{X})).$$

We want to have this Loss as small as possible. Well, under our setting, this will be a R.V., so we define the **Average Loss** or the **Risk** of the Predictor $g$ to be

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

Here the Expectation is over the Distribution of $(\mathbf{X}, Y)$, i.e., $\mathcal{F}$.

Now, we can state our Problem of finding a good Predictor: Find $g$ minimizing the Risk, i.e., find

$$g^* \in \underset{g}{argmin}\, Risk(g).$$

## Example

**Toy Example:** Assume $\mathcal{X} = \{1, 2, 3\}$, $\mathcal{Y} = \{0, 1\}$, and we have the Joint Distribution of $(X, Y)$:

| $Y \setminus X$ | 1 | 2 | 3 |
|:---:|:---:|:---:|:---:|
| 0 | 0.1 | 0.2 | 0.1 |
| 1 | 0.2 | 0.1 | 0.3 |

Assume

$$g(x) = \begin{cases} 0, & \text{if } x \text{ is even} \\ 1, & \text{otherwise} \end{cases}$$

and $\ell(y_1, y_2) = |y_1 - y_2|$. Calculate the $Risk(g)$.

**Solution:** OTB

# Predictors

The above Minimization Problem is not complete: we need to specify the set of all possible Predictors $g$.

# Predictors

The above Minimization Problem is not complete: we need to specify the set of all possible Predictors $g$.

The best case will be if we will take $g$ to be **any measurable function from $\mathcal{X}$ to $\mathcal{Y}$**.

# Predictors

The above Minimization Problem is not complete: we need to specify the set of all possible Predictors $g$.

The best case will be if we will take $g$ to be **any measurable function from $\mathcal{X}$ to $\mathcal{Y}$**. But, unfortunately, this set is veery large to be able to solve the problem there.

# Predictors

The above Minimization Problem is not complete: we need to specify the set of all possible Predictors $g$.

The best case will be if we will take $g$ to be **any measurable function from $\mathcal{X}$ to $\mathcal{Y}$**. But, unfortunately, this set is veery large to be able to solve the problem there.

Usually, we assume that $g$ comes from a Parametric Family of functions, which we call a Predictive Model:

$$g \in \mathcal{G} = \{g(\mathbf{x}|\theta),\ \theta \in \Theta\}, \qquad where \quad g(\mathbf{x}|\theta) : \mathcal{X} \to \mathcal{Y}.$$

and $\Theta$ is some Parameter Set (1D or more).

# Predictors, Examples

Say,

- In the 1D Linear Regression Problem, we consider

$$g(\mathbf{x}|\theta) = \theta_0 + \theta^T \cdot \mathbf{x}$$

# Predictors, Examples

Say,

- In the 1D Linear Regression Problem, we consider

$$g(\mathbf{x}|\theta) = \theta_0 + \theta^T \cdot \mathbf{x} = \theta_0 + \theta_1 \cdot x_1 + ... + \theta_d \cdot x_d;$$

# Predictors, Examples

Say,

- In the 1D Linear Regression Problem, we consider

$$g(\mathbf{x}|\theta) = \theta_0 + \theta^T \cdot \mathbf{x} = \theta_0 + \theta_1 \cdot x_1 + ... + \theta_d \cdot x_d;$$

- In the Binary Classification Problem, with $\mathcal{Y} = \{-1, 1\}$, we can consider, say

$$g(\mathbf{x}|\theta) = sgn(\theta_0 + \theta_1 \cdot x_1 + ... + \theta_d \cdot x_d);$$

# Predictors, Examples

Say,

- In the 1D Linear Regression Problem, we consider

$$g(\mathbf{x}|\theta) = \theta_0 + \theta^T \cdot \mathbf{x} = \theta_0 + \theta_1 \cdot x_1 + ... + \theta_d \cdot x_d;$$

- In the Binary Classification Problem, with $\mathcal{Y} = \{-1, 1\}$, we can consider, say

$$g(\mathbf{x}|\theta) = sgn(\theta_0 + \theta_1 \cdot x_1 + ... + \theta_d \cdot x_d);$$

- In the general Regression/Classification Problems, we can have $g(\mathbf{x}|\theta)$ to be a Neural Network, where $\mathbf{x}$ is our input, $\theta$ is the vector of all NN weights, and $g(\mathbf{x}|\theta)$ is the output of the NN.

# The Learning Problem

Now we can finalize the statement of our Problem:

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- ▶ A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- ▶ A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- ▶ A Predictive Model (set of Functions) $\mathcal{G}$

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- A Predictive Model (set of Functions) $\mathcal{G}$

and we want to find $g^* \in \mathcal{G}$ such that

$$g^* \in \underset{g \in \mathcal{G}}{\text{argmin}} \, Risk(g) = \underset{g \in \mathcal{G}}{\text{argmin}} \, \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- A Predictive Model (set of Functions) $\mathcal{G}$

and we want to find $g^* \in \mathcal{G}$ such that

$$g^* \in \underset{g \in \mathcal{G}}{\arg\min}\, Risk(g) = \underset{g \in \mathcal{G}}{\arg\min}\, \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

If $\mathcal{G}$ coincides with the set of all measurable functions, thet $g^*$, if exists, is called the **Bayes Predictor**.

# The Learning Problem

Now we can finalize the statement of our Problem: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- A Predictive Model (set of Functions) $\mathcal{G}$

and we want to find $g^* \in \mathcal{G}$ such that

$$g^* \in \underset{g \in \mathcal{G}}{\text{argmin}} \, Risk(g) = \underset{g \in \mathcal{G}}{\text{argmin}} \, \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big).$$

If $\mathcal{G}$ coincides with the set of all measurable functions, thet $g^*$, if exists, is called the **Bayes Predictor**. And if $g^*$ is a Bayes Predictor, then its Risk, $Risk(g^*)$, is called the **Bayes Risk**.

# Risk Minimization vs Empirical Risk Minimization

OK, nice. But, unfortunately, we usually cannot solve the above problem, since we **do not have the Distribution** $\mathcal{F}$ to calculate the Risk.

# Risk Minimization vs Empirical Risk Minimization

OK, nice. But, unfortunately, we usually cannot solve the above problem, since we **do not have the Distribution** $\mathcal{F}$ to calculate the Risk.

So we do the following: recall that, by the LLN,

$$\frac{1}{n} \cdot \sum_{k=1}^{n} \ell(Y_k, g(\mathbf{X}_k)) \rightarrow$$

# Risk Minimization vs Empirical Risk Minimization

OK, nice. But, unfortunately, we usually cannot solve the above problem, since we **do not have the Distribution** $\mathcal{F}$ to calculate the Risk.

So we do the following: recall that, by the LLN,

$$\frac{1}{n} \cdot \sum_{k=1}^{n} \ell(Y_k, g(\mathbf{X}_k)) \to \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = Risk(g) \qquad a.s.$$

# Risk Minimization vs Empirical Risk Minimization

OK, nice. But, unfortunately, we usually cannot solve the above problem, since we **do not have the Distribution** $\mathcal{F}$ to calculate the Risk.

So we do the following: recall that, by the LLN,

$$\frac{1}{n} \cdot \sum_{k=1}^{n} \ell(Y_k, g(\mathbf{X}_k)) \to \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = Risk(g) \qquad a.s.$$

So, instead of trying to minimize $Risk(g)$, we can try to minimize

$$ERM(g) = \frac{1}{n} \cdot \sum_{k=1}^{n} \ell(Y_k, g(\mathbf{X}_k)),$$

which is called the **Empirical Risk Measure of** $g$.

Now we change the statement of our Problem like this:

# The Learning Problem, Empirical Version

Now we change the statement of our Problem like this: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- A Predictive Model (set of Functions) $\mathcal{G}$

# The Learning Problem, Empirical Version

Now we change the statement of our Problem like this: We are given

- A Dataset of Observations $(\mathbf{x}_k, y_k)$, $k = 1, .., n$, coming as a realization of $(\mathbf{X}_k, Y_k)$ from an unknown Distribution $\mathcal{F}$;

- A Loss Function $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$;

- A Predictive Model (set of Functions) $\mathcal{G}$

and we want to find $g^* \in \mathcal{G}$ such that

$$g^* \in \underset{g \in \mathcal{G}}{argmin}\, ERM(g).$$

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);
- our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);
- our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

Now, if $g$ is any Predictor (any function $g : \mathcal{X} \to \mathcal{Y}$), then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big(\mathbf{1}(Y \neq g(\mathbf{X}))\Big) =$$

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);
- our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

Now, if $g$ is any Predictor (any function $g : \mathcal{X} \to \mathcal{Y}$), then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big(\mathbf{1}(Y \neq g(\mathbf{X}))\Big) = \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);
- our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

Now, if $g$ is any Predictor (any function $g : \mathcal{X} \to \mathcal{Y}$), then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big(\mathbf{1}(Y \neq g(\mathbf{X}))\Big) = \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

So the Risk of $g$ in the case of Binary Classification with $0 - 1$ Loss is **the Probability to predict incorrectly**.

## Example: Binary Classification Problem

In this case we have:

- ▶ $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);

- ▶ our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

Now, if $g$ is any Predictor (any function $g : \mathcal{X} \to \mathcal{Y}$), then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big(\mathbf{1}(Y \neq g(\mathbf{X}))\Big) = \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

So the Risk of $g$ in the case of Binary Classification with $0 - 1$ Loss is **the Probability to predict incorrectly**.

And the Problem to find the Bayes Predictor is to find a function $g$ with minimal Probability of incorrect Prediction:

$$g^* \in \underset{g}{argmin}\, \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

# Example: Binary Classification Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} = \{0, 1\}$ (or, $\mathcal{Y} = \{-1, 1\}$);

- our Loss Function is the $0 - 1$ loss: $\ell(y_1, y_2) = \mathbf{1}(y_1 \neq y_2)$

Now, if $g$ is any Predictor (any function $g : \mathcal{X} \to \mathcal{Y}$), then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big(\mathbf{1}(Y \neq g(\mathbf{X}))\Big) = \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

So the Risk of $g$ in the case of Binary Classification with $0 - 1$ Loss is **the Probability to predict incorrectly**.

And the Problem to find the Bayes Predictor is to find a function $g$ with minimal Probability of incorrect Prediction:

$$g^* \in \underset{g}{argmin} \, \mathbb{P}\Big(Y \neq g(\mathbf{X})\Big).$$

Can you guess $g^*$ ? ☺

# Example: Least Squares Regression Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} \subset \mathbb{R}$;

# Example: Least Squares Regression Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} \subset \mathbb{R}$;
- our Loss Function is the Quadratic Loss: $\ell(y_1, y_2) = (y_1 - y_2)^2$

# Example: Least Squares Regression Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} \subset \mathbb{R}$;
- our Loss Function is the Quadratic Loss: $\ell(y_1, y_2) = (y_1 - y_2)^2$

And, if $g$ is any Predictor, then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big((Y - g(\mathbf{X}))^2\Big).$$

# Example: Least Squares Regression Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} \subset \mathbb{R}$;
- our Loss Function is the Quadratic Loss: $\ell(y_1, y_2) = (y_1 - y_2)^2$

And, if $g$ is any Predictor, then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big((Y - g(\mathbf{X}))^2\Big).$$

The Problem to find the Bayes Predictor in this case is:

$$g^* \in \underset{g}{argmin}\, \mathbb{E}\Big((Y - g(\mathbf{X}))^2\Big).$$

# Example: Least Squares Regression Problem

In this case we have:

- $\mathcal{X}$ is arbitrary, $\mathcal{Y} \subset \mathbb{R}$;
- our Loss Function is the Quadratic Loss: $\ell(y_1, y_2) = (y_1 - y_2)^2$

And, if $g$ is any Predictor, then

$$Risk(g) = \mathbb{E}\Big(\ell(Y, g(\mathbf{X}))\Big) = \mathbb{E}\Big((Y - g(\mathbf{X}))^2\Big).$$

The Problem to find the Bayes Predictor in this case is:

$$g^* \in \underset{g}{argmin}\, \mathbb{E}\Big((Y - g(\mathbf{X}))^2\Big).$$

Can you guess $g^*$ ? ⌣

# Conditional Distribution, Discrete Case

To continue with finding Bayes Predictors, we need the notion of Conditional Distribution.

# Conditional Distribution, Discrete Case

To continue with finding Bayes Predictors, we need the notion of Conditional Distribution.

First assume $X$ and $Y$ are Jointly Distributed Discrete r.v. with the PMF $f_{X,Y}(x,y)$, i.e.,

$$f_{X,Y}(x,y) = \mathbb{P}(X = x, Y = y).$$

# Conditional Distribution, Discrete Case

To continue with finding Bayes Predictors, we need the notion of Conditional Distribution.

First assume $X$ and $Y$ are Jointly Distributed Discrete r.v. with the PMF $f_{X,Y}(x,y)$, i.e.,

$$f_{X,Y}(x,y) = \mathbb{P}(X = x, Y = y).$$

The Marginal PMFs for $X$ and $Y$ will be:

$$f_X(x) = \mathbb{P}(X = x) = \sum_y f_{X,Y}(x,y),$$

$$f_Y(y) = \mathbb{P}(Y = y) = \sum_x f_{X,Y}(x,y).$$

## Conditional Distribution, Discrete Case

To continue with finding Bayes Predictors, we need the notion of Conditional Distribution.

First assume $X$ and $Y$ are Jointly Distributed Discrete r.v. with the PMF $f_{X,Y}(x, y)$, i.e.,

$$f_{X,Y}(x, y) = \mathbb{P}(X = x, Y = y).$$

The Marginal PMFs for $X$ and $Y$ will be:

$$f_X(x) = \mathbb{P}(X = x) = \sum_y f_{X,Y}(x, y),$$

$$f_Y(y) = \mathbb{P}(Y = y) = \sum_x f_{X,Y}(x, y).$$

Next, we know that

$$\mathbb{P}(Y = y | X = x) =$$

# Conditional Distribution, Discrete Case

To continue with finding Bayes Predictors, we need the notion of Conditional Distribution.

First assume $X$ and $Y$ are Jointly Distributed Discrete r.v. with the PMF $f_{X,Y}(x, y)$, i.e.,

$$f_{X,Y}(x, y) = \mathbb{P}(X = x, Y = y).$$

The Marginal PMFs for $X$ and $Y$ will be:

$$f_X(x) = \mathbb{P}(X = x) = \sum_y f_{X,Y}(x, y),$$

$$f_Y(y) = \mathbb{P}(Y = y) = \sum_x f_{X,Y}(x, y).$$

Next, we know that

$$\mathbb{P}(Y = y | X = x) = \frac{\mathbb{P}(Y = y, X = x)}{\mathbb{P}(X = x)},$$

# Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

## Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

## Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) =$$

## Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

Now, we can calculate, say

$$\mathbb{P}(Y \in A|X = x) = \sum_{y \in A} f_{Y|X}(y|x) =$$

# Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \sum_{y \in A} f_{Y|X}(y|x) = \sum_{y \in A} f_{Y|X=x}(y).$$

## Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

Now, we can calculate, say

$$\mathbb{P}(Y \in A|X = x) = \sum_{y \in A} f_{Y|X}(y|x) = \sum_{y \in A} f_{Y|X=x}(y).$$

Also, we can calculate, say, the Expected value of $Y$ given the value of $X$:

$$\mathbb{E}(Y|X = x) =$$

## Conditional Distribution, Discrete Case

Hence, we define the Conditional PMF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\sum_y f_{X,Y}(x,y)}.$$

Sometimes we write this as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \sum_{y \in A} f_{Y|X}(y|x) = \sum_{y \in A} f_{Y|X=x}(y).$$

Also, we can calculate, say, the Expected value of $Y$ given the value of $X$:

$$\mathbb{E}(Y|X = x) = \sum_y y \cdot f_{Y|X}(y|x) = \sum_y y \cdot f_{Y|X=x}(y)$$

# Example

# Conditional Distribution, Continuous Case

Now, if $X$ and $Y$ are Jointly Distributed Continuous r.v. with the PDF $f_{X,Y}(x,y)$, i.e.,

$$\mathbb{P}((X,Y) \in A) = \iint_{(x,y) \in A} f_{X,Y}(x,y)dxdy.$$

# Conditional Distribution, Continuous Case

Now, if $X$ and $Y$ are Jointly Distributed Continuous r.v. with the PDF $f_{X,Y}(x, y)$, i.e.,

$$\mathbb{P}((X, Y) \in A) = \iint_{(x,y) \in A} f_{X,Y}(x, y) dx dy.$$

Then the Marginal PDFs for $X$ and $Y$ will be:

$$f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x, y) dy, \qquad f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x, y) dx.$$

## Conditional Distribution, Continuous Case

Now, if $X$ and $Y$ are Jointly Distributed Continuous r.v. with the PDF $f_{X,Y}(x,y)$, i.e.,

$$\mathbb{P}((X,Y) \in A) = \iint_{(x,y) \in A} f_{X,Y}(x,y)dxdy.$$

Then the Marginal PDFs for $X$ and $Y$ will be:

$$f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x,y)dy, \qquad f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y)dx.$$

Now, in the analogy of the Discrete case, we define the Conditional PDF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\int_{\mathbb{R}} f_{X,Y}(x,y)dy},$$

for all $x$ such that $f_X(x) \neq 0$.

# Conditional Distribution, Continuous Case

Now, if $X$ and $Y$ are Jointly Distributed Continuous r.v. with the PDF $f_{X,Y}(x,y)$, i.e.,

$$\mathbb{P}((X,Y) \in A) = \iint_{(x,y) \in A} f_{X,Y}(x,y) dx dy.$$

Then the Marginal PDFs for $X$ and $Y$ will be:

$$f_X(x) = \int_{\mathbb{R}} f_{X,Y}(x,y) dy, \qquad f_Y(y) = \int_{\mathbb{R}} f_{X,Y}(x,y) dx.$$

Now, in the analogy of the Discrete case, we define the Conditional PDF of $Y$ given $X$ by

$$f_{Y|X}(y|x) = \frac{f_{X,Y}(x,y)}{f_X(x)} = \frac{f_{X,Y}(x,y)}{\int_{\mathbb{R}} f_{X,Y}(x,y) dy},$$

for all $x$ such that $f_X(x) \neq 0$. Again, we write this also as

$$f_{Y|X=x}(y) = \frac{f_{X,Y}(x,y)}{f_X(x)}.$$

# Conditional Distribution, Continuous Case

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) =$$

# Conditional Distribution, Continuous Case

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \int_A f_{Y|X}(y|x)\,dy =$$

# Conditional Distribution, Continuous Case

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \int_A f_{Y|X}(y|x) \, dy = \int_A f_{Y|X=x}(y) \, dy.$$

# Conditional Distribution, Continuous Case

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \int_A f_{Y|X}(y|x)\, dy = \int_A f_{Y|X=x}(y)\, dy.$$

And the Expected value of $Y$ given the value of $X$ will be:

$$\mathbb{E}(Y|X = x) =$$

# Conditional Distribution, Continuous Case

Now, we can calculate, say

$$\mathbb{P}(Y \in A | X = x) = \int_A f_{Y|X}(y|x)\,dy = \int_A f_{Y|X=x}(y)\,dy.$$

And the Expected value of $Y$ given the value of $X$ will be:

$$\mathbb{E}(Y|X = x) = \int_{\mathbb{R}} y \cdot f_{Y|X}(y|x)\,dy = \int_{\mathbb{R}} y \cdot f_{Y|X=x}(y)\,dy.$$

# Example

**Example:** Say, $(X, Y) \sim \text{Unif}([0, 1] \times [0, 2])$. What is $Y | X = 1$, or, in general, $Y | X = x$?
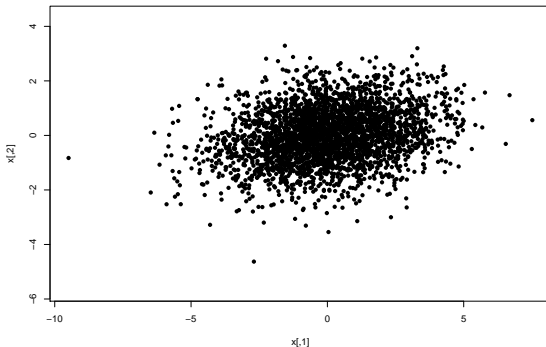
# Example

**Example:** Say, $(X, Y) \sim Unif([0, 1] \times [0, 2])$. What is $Y|X = 1$, or, in general, $Y|X = x$?

**Example:** Now, assume $(X, Y) \sim Unif(D)$, where $D$ is the triangle with vertices at $(-1, 0)$, $(0, 1)$ and $(1, 0)$. What is $Y|X = x$ ?

## Example

**Example:** Assume $(X, Y) \sim \mathcal{N}(\mu, \Sigma)$, where

$$\mu = \left[ \begin{array}{c} 0 \\ 0 \end{array} \right], \qquad \Sigma = \left[ \begin{array}{cc} 4 & 0.5 \\ 0.5 & 1 \end{array} \right]$$



What is $Y|X = x$?