# Information Theory Recitation

Siming He, ESE 5460

November 2023

## Contents

## 1  Elements of Information Theory

Information theory was developed as a mathematical theory of communication. A communication system can be divided into five parts:

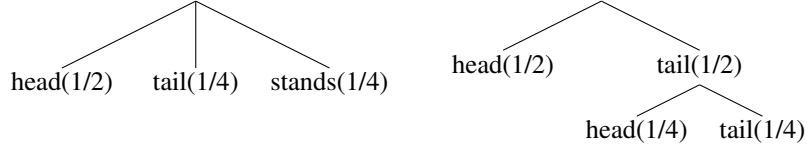**source → transmitter → channel → receiver → destination**

1. **Source** is the message that we want to send to the destination, e.g. text messages, images, videos

2. **Transmitter** converts the message into some signals that can be transmitted. For example, Morse code converts numbers and English alphabets into dots and dashes. Dots and dashes signals can be easily transmitted by short-duration light/sound and long-duration light/sound.

3. **Channel** is the medium used to transmit signal. For example, air transmits light/sound. And wires transmits pulses of voltages.

4. **Receiver** performs an inverse operation of transmitter to reconstruct messages from signals.

5. **Destination** would be the person who receives and reads your text.

Two fundamental questions that are essential to communication system and can be answered by information theory.

1. **What is the maximal compression of data?** We want the communication system to be energy-efficient, i.e., spend the minimum amount of energy to transmit messages. To achieve this, we want to (maximally and losslessly) **compress** our messages.

2. **What is the maximum rate that data can be transmitted?** The maximally compressed message would work well if the channel is noiseless. However, if the channel is noisy, compressed messages may be easily corrupted by the noise during transmission. To make the signal robust to noise, we also want to add **redundancy** into the signal. So the maximum rate of transmission depends on both compression and necessary redundancy.

To understand the design of transmitter or encoder, we have to understand the information in the source or input. However, the concept of information is vague and broad and is hard to be formulated mathematically. Instead, we can define a measure of the amount of information in a random variable. For a probability distribution $P^{(n)} = (p_1, \cdots, p_n)$ over an alphabet of size $n$, we want to find an information measure $H_n(p_1, \cdots, p_n)$ characterising the information content of chance experiment with $P^{(n)}$. Shannon argued that there are three required properties of desirable information measure:

1. $H_n(p_1, \cdots, p_n)$ is a continuous function of $p_1, \cdots, p_n$. This is intuitive since a small change in the probability distribution shouldn't lead to a sudden jump in information content.

2. For $H_n(p_1, \cdots, p_n)$ where $p_1 = \cdots = p_n$, $H_n$ should be a monotonic increasing function of $n$. Intuitively, more events leads to more uncertainty and more information. Rolling a dice seems to be more complex than flipping a coin.

3. The measure should have $H_n(p_1, \cdots, p_n) = H_{n-1}(p_1 + p_2, p_3, \cdots, p_n) + (p_1 + p_2)H_2(\frac{p_1}{p_1+p_2}, \frac{p_2}{p_1+p_2})$. For example, $H_3(\frac{1}{2}, \frac{1}{4}, \frac{1}{4})$ should be equal to $H_2(\frac{1}{2}, \frac{1}{2}) + \frac{1}{2}H_2(\frac{1}{2}, \frac{1}{2})$. Intuitively, if the information measure measures the amount of information in a chance outcome, how the source produces such outcome shouldn't change the amount of information. For example, the source can flip a magical coin that is head with probability $\frac{1}{2}$, is tail with probability $\frac{1}{4}$, stands up with probability $\frac{1}{4}$. The source can also flip a fair coin and flip the coin again if the first flip results in tail. Those two process gives the same chance outcome so should contain the same amount of information. For the second process, we firstly gain the information of $H_2(\frac{1}{2}, \frac{1}{2})$. Then with probability $\frac{1}{2}$, we gain additional information $H_2(\frac{1}{2}, \frac{1}{2})$ from the second round.

head(1/2)    tail(1/4)    stands(1/4)        head(1/2)              tail(1/2)

head(1/4)    tail(1/4)

**Definition 1.1** (Entropy)**.** The only $H$ satisfying the three properties above is in the form of $H = -K\sum_{i=1}^{n} p_i \log p_i$. $K$ is basically the unit of the measure and can be 1. Then, we have the definition of entropy:

$$H_n(p_1, \cdots, p_n) = -\sum_{i=1}^{n} p_i \log p_i$$

For a random variable $X$, we can also write the entropy as

$$H(X) = -\sum_{x} p(x) \log p(x)$$

*Remark.* The quantity $-\log p(x)$ can be viewed as the information content of event $x$. The quantity aligns with our intuition that rare events should contain more information $-\log p(x) > -\log p(y)$ when $p(x) < p(y)$. Entropy can be viewed as the expected information content from a random variable.

*Remark.* Another view of entropy is that entropy of a random variable is a lower bound on the average length of the shortest description of the random variable. For each event $x$, we can encode it by $l_x$ bits. Then, the expected length of the description of $X$ is $\sum_x p(x)l_x$. Gibbs' inequality shows that when $l_x = -\log_2 p(x)$, we have the shortest expected length. Specifically, the inequality states $-\sum_x p(x) \log p(x) \leq -\sum_x p(x) \log q(x)$ for any description with length $l_x = -\log q(x)$. It's reasonable to use the shortest expected length of description as a measure of information, since a process with more information inevitably requires more bits to describe.

**Definition 1.2** (Joint Entropy and Conditional Entropy)**.** The joint entropy of a pair of discrete variables $(X, Y)$ with joint distribution $p(x, y)$ is defined as

$$H(X, Y) = -\sum_{x} \sum_{y} p(x, y) \log p(x, y)$$

The conditional entropy $H(y|X)$ is then defined as

$$H(Y|X) = -\sum_{x} p(x) \sum_{y} p(y|x) \log p(y|x)$$
$$= \sum_{x} p(x) H(Y|X = x)$$

Following the chain rule of probability, we also have the chain rule for entropy:

**Theorem 1** (Chain Rule).

$$H(X, Y) = H(X) + H(Y|X)$$

**Definition 1.3** (Relative Entropy / Kullback–Leibler(KL) Divergence). The relative entropy between two probability mass functions $p(x)$ and $q(x)$ is

$$D(p||q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$$

KL Divergence measures the dissimilarity between two probability mass functions. However, KL Divergence is not a distance measure since it's not symmetric.
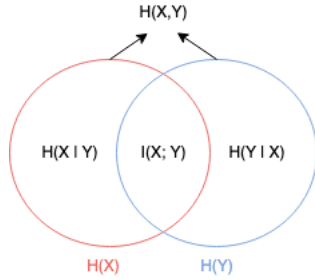
**Theorem 2.**

$$D(p||q) \geq 0$$

*and*

$$D(p||q) = 0 \ \textit{if and only if} \ p = q$$

**Definition 1.4** (Mutual Information). The mutual information between two random variables $X, Y$ is

$$
\begin{aligned}
I(X; Y) &= \sum_x \sum_y p(x, y) \log \frac{p(x, y)}{p(x)p(y)} \\
&= D(p(x, y)||p(x)p(y)) \\
&= H(X) + H(Y) - H(X, Y) \\
&= H(X) - H(X|Y) \\
&= H(Y) - H(Y|X)
\end{aligned}
$$

Mutual information measures the amount of information that one variable contains about another. This measure is symmetric.



This image demonstrates the relationship between mutual information and entropy. $H(X, Y)$ can be viewed as the information contained in joint $(X, Y)$. $H(X, Y) \leq H(X) + H(Y)$ since $X$ and $Y$ may have some shared information $I(X, Y)$. $H(X|Y)$ measures the amount of additional information we get from $X$ if we already know $Y$, so we need to remove the amount of shared information: $H(X|Y) = H(X) - I(X; Y)$.

**Theorem 3.**

$$I(X;Y) \geq 0$$

*and*

$$I(X;Y) = 0 \ if \ and \ only \ if \ X \ and \ Y \ are \ independent$$

**Theorem 4** (Data-processing Inequality). *If we have of Markov chain of three random variables $X \to Y \to Z$, i.e., $Z$ is conditionally independent from $X$ given $Y$, then*

$$I(X;Y) \geq I(X;Z)$$

*Remark.* This is a really cool theorem! There is no data processing $Y \to Z$ such that we can get more information about $X$. Consider image classification as an example. Let $X$ be the label of images. Let $Y$ be images generated from the labels. And let $f$ be a neural network converting $Y$ into some embedding $Z = f(Y)$. We have a Markov chain $X \to Y \to f(Y)$. And based on this theorem, we know that $I(X;Y) \geq I(X;f(Y))$. It means that neural networks are not processing images to gain more information about $Y$! So what are neural network actually doing from an information-theoretic view? We will soon talk about information bottleneck which answers the question to some extent.

**Definition 1.5** (Sufficient Statistics). Given a family of probability mass functions $\{f_\theta(x)\}$. A statistic $T(X)$ is a function of the sample $X$ from $\{f_\theta(x)\}$, e.g., sample mean or sample variance. We have a Markov chain $\theta \to X \to T(X)$. By Data-processing Inequality, we have

$$I(\theta, T(X)) \leq I(\theta, X)$$

$T(X)$ is called sufficient statistic if it contains all the information that $X$ has about $\theta$, i.e.,

$$I(\theta, T(X)) = I(\theta, X)$$

*Remark.* Taking the same example of classification. If the neural network $f$ is a sufficient statistic, then we would have $I(X;Y) = I(X;f(Y))$. It means the embedding $f(Y)$ and the original image $Y$ contain the same amount of information about label $X$.

## 2 Information Bottleneck in Deep Learning

In a typical supervised learning problem, we use input data from $X$ to predict output label from $Y$. Ideally, we want the neural network to learn some embedding $\hat{X}$ such that all the information relevant to $Y$ is kept and other information is discarded. We consider the Markov chain $Y \to X \to \hat{X}$. To discard
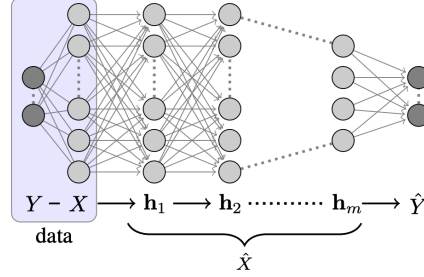
Figure 1: Source: Deep Learning and the Information Bottleneck Principle

irrelevant information, we want to maximally compress $X$ so that only information about $Y$ is kept. However, for a lossy compression, by Data-processing Inequality, $I(Y, \hat{X}) \leq I(Y, X)$. It means that $\hat{X}$ will loss information about $Y$. Obviously, there is a trade-off between compressing the representation and preserving meaningful information about $Y$. Information bottleneck method is a principle designed to extract relevant information in $X$ about $Y$. It's called "information bottleneck" because we are passing the information in $X$ about $Y$ through a "bottleneck" (e.g., a neural network) to form a more compact representation $\hat{X}$.

**Definition 2.1** (Information Bottleneck Principle)**.** To find a optimal representation $\hat{X}$, we want to minimize the functional

$$I(\hat{X}, X) - \beta I(\hat{X}, Y)$$

*Remark.*    1. Minimizing $I(\hat{X}, X)$ is maximizing compression. Without compression ($\hat{X}$ has the same information as $X$), we have $I(\hat{X}, X) = I(X, X) = H(X)$. With maximal compression ($\hat{X}$ only contains information of $Y$), we have $I(\hat{X}, X) = I(Y, X)$. And $I(Y, X) = H(X) - H(X|Y) \leq H(X)$.

2. Maximizing $I(\hat{X}, Y)$ is maximizing the relevant information in $\hat{X}$ about $Y$. This is obvious from the definition of mutual information.

3. $\beta$ is the Lagrange multiplier. If $\beta = 0$, $X$ will be maximally compressed at the risk of losing a lot of relevant information about $Y$. If $\beta \to \infty$, $\hat{X}$ would be a super detailed representation of $X$ without any compression.

Now, if we consider an actual neural network as follow: This network forms the Markov chain $Y \to X \to h_1 \to \cdots \to h_m \to \hat{Y}$. By Data-processing inequality, we have

$$H(X) = I(X, X) \geq I(X, h_1) \geq \cdots \geq I(X, h_m) \geq I(X, \hat{Y})$$
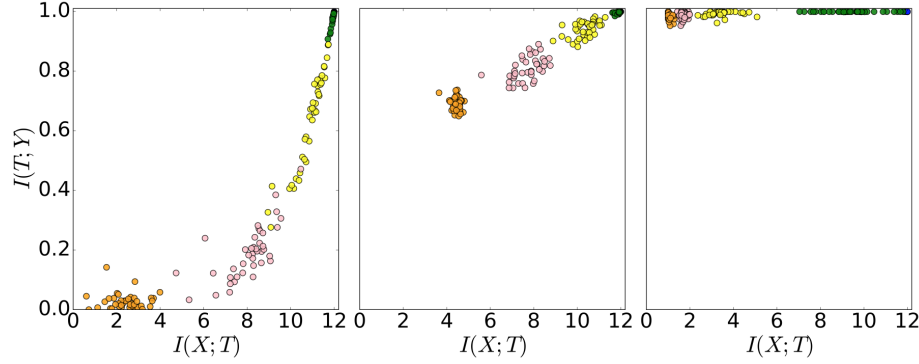
and

$$I(Y, X) \geq I(Y, h_1) \geq \cdots \geq I(Y, h_m) \geq I(Y, \hat{Y})$$

By Information Bottleneck Principle, we want to maximize $I(Y, \hat{Y})$ while minimizing $I(X, \hat{Y})$. For each layer, we want to maximize $I(Y, h_i)$ while minimizing $I(h_{i-1}, h_i)$.

**Theorem 5** (Information Plane). *Most supervise learning algorithms try to balance the trade-off between the performance on sampled data and the generalization to unseen data. Naftali Tishby argues that it's reasonable to consider $I(Y, T)$ as a measure of performance and consider $I(X, T)$ as a control of generalization.*

## 2.1 Optimization Phases of Neural Network



Source: OPENING THE BLACK BOX OF DEEP NEURAL NETWORKS VIA INFORMATION.

## 2.2 Input Compression Bound

**Definition 2.2** (Vapnik-Chernovenkis (VC) bound). As learned in lecture, the bound is

$$R(f_{ERM}) \le R(f^*) + 2\sqrt{\frac{1}{2n} \log \frac{4|\mathcal{F}|}{\delta}}$$

where $n$ is the number of samples, $\delta$ is the confidence, and $|\mathcal{F}|$ is the size of the hypothesis class. This bound becomes vacuous for neural networks which have super large $|\mathcal{F}|$.

**Definition 2.3** (Input Compression Bound). We want to get a bound that depends on the information of input instead of the large hypothesis class. The bound we will get is

$$R(f_{ERM}) \le R(f^*) + 2\sqrt{\frac{1}{2n} \left( 2^{I(T_\epsilon, X)} + \log \frac{1}{\delta} \right)}$$

7

where $T_\epsilon$ is a $\epsilon$ partition of $X$. $T_\epsilon$ can be considered as embedding space where each $X$ is compressed into some code. Note that if we can reduce $I(T_\epsilon, X)$ by $k$ ($k$ more bits of compression), we can reduce the training samples by a factor of $2^k$!



This figure above shows the partition of the input space. We partition the input space $X$ into small red balls $b_1, b_2, \cdots$. Each red ball can either be 0 or 1. Under this partition, each $x \in X$ has a more compact binary representation, e.g., $b_1 = 0, b_2 = 1, b_3 = 1, \cdots$. We can think of supervise learning as creating a good partition of $X$ that aligns with the labels $Y$. The typical set of $X$ has size $2^{H(X)}$ (you can read more in Elements of Information Theory 3.29). Similarly, each red ball has size $2^{H(X|T_\epsilon)}$. Then the size of $T_\epsilon$ (the number of red balls) is $|T_\epsilon| = \frac{2^{H(X)}}{2^{H(X|T_\epsilon)}} = 2^{I(X,T_\epsilon)}$. If we consider the hypothesis class of Boolean function that maps $x$ to the binary representation $b_1 = 0, b_2 = 1, \cdots, b_{|T_\epsilon|}$, this class has size $2^{2^{I(X,T_\epsilon)}}$. Then, we can get the input compression bound from $\log 2^{2^{I(X,T_\epsilon)}}$.

## 2.3  $I(T, Y)$ and Performance



Source

If we know $X, Y$, we will have optimal lossy compression which is the black line. Then, the model's performance is subject to the compression loss. Since we only have finite data, the compression is worse and has the red line as the lower bound. The model's performance is subject to an additional finite sample loss.

# 3  References

1. A Mathematical Theory of Communication

2. Elements of Information Theory

3. Deep Learning and the Information Bottleneck Principle

4. Opening the black box of Deep Neural Networks via Information

5. Learning and generalization with the information bottleneck

6. A great unpublished note from Santosh Venkatesh on An Axiomatic Approach to Information

# AWS Tutorial

ESE 5460

October 2023

## Contents

## 1 Overview

The EC2 (Elastic Cloud Compute) service on AWS (https://aws.amazon.com) offers cloud computing infrastructure and provides access to GPU resources. Each student will be given compute credits to use on AWS. AWS could be used for the programming assignments and is a useful tool for the projects. EC2 gives you access to GPU instances which are charged by the hour. In this short tutorial, we will go over the basics, best practices and some useful tools

# 2 Why use the Cloud?

- Use latest GPU architectures
- Easier to scale
- No maintenance costs

# 3    Creating an Account

You will get an email titled "AWS @ Penn access" from isc-cloud-solutions@isc.upenn.edu with instructions on how to create your account and use the AWS credits allocated to ESE 546.

# 4    Launching a GPU instance

## 4.1    Use Penn Campus Network or University VPN

VPN Link

## 4.2    Switch your zone to Virginia (U.S East) on the top right of the screen.

We are going to use this zone and it generally gives you better availability for the machines.



## 4.3    Name and Tags

Search and go to the EC2 service in AWS and click on **Launch instance** to begin.

**In Name and tags**,

1. Give a name to your instance (it can be anything but it would be good to include your name e.g. pratikac-instance)

2. Click on **Additional Tags**, create a key called "Owner" and value as "your Penn email including everything after @" (i.e., pratikac@seas.upenn.edu for Pratik). Then add Intsances and Volumes to the Resource types. *We have scripts to prevent people from deleting each other's instances accidentally. So you will not be able to interact with any instance or volume that was not launched with "Owner:your-key".*

## 4.4   Select Amazon Machine Image (AMI)

Under **Application and OS Images**, click **Browse more AMIs**, search **Deep Learning AMI**, and select "Deep Learning AMI GPU PyTorch ... (Ubuntu ...)".



You can also use a standard Ubuntu but you will have to install everything (deep learning environment and libraries) yourself. so unless you know what you are doing, I would advise using the AMI.
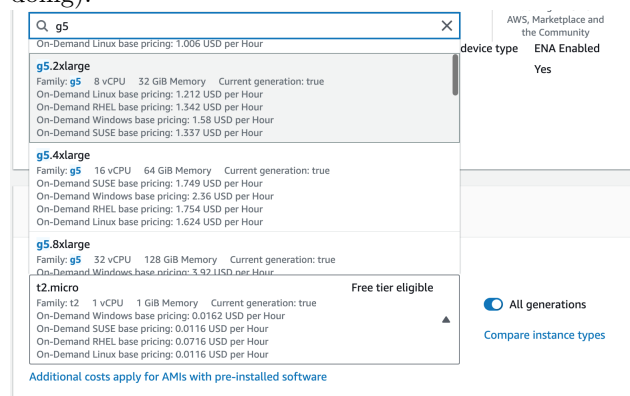
## 4.5 Choose Instance Type

Under **Instance type**, change the instance type to one of the GPU instance. **You should use the following instances**:

1. g5.xlarge: 1 Nvidia A30 GPU (4 CPU cores, 16 GB memory). This should be sufficient for almost everyone and it is also the cheapest.

2. g5.2xlarge: 1 Nvidia A30 GPU (8 CPU cores, 32 GB memory). If you want run a dataloader with many threads, e.g., for lots of augmentation.

3. g5.12xlarge: 4 Nvidia A30 GPUs (48 CPU cores, 192 GB memory). Do not use this in general. Use it only if you need all 4 GPUs. This instance costs  12 dollars/hour so you will use your 100-150 dollars or so credits very quickly.

4. You can also use g4 instances and they are only a tiny bit slower (these have Nvidia T4 GPU which is the same one that Google Colab, so they are equally fast).

You can run multiple training runs simultaneously on the same GPU without changing any of your code. This is not different from running multiple processes on the same CPU. You do not need multiple GPUs just because you are training multiple neural networks. **Do not use the other g5 instances** unless you know for sure that you will benefit from them (and you know what you are doing).



## 4.6 Create a keypair from the EC2 console

Under the Key pair (login) section, **Create new key pair**. Without this keypair you will not be able to SSH into the instance and will not be able to use it. It is a good idea to call the keypair by your Penn email, e.g., pratikac for Pratik.

## 4.7 Security Groups

Under **Network Settings**, choose the option to **Select existing security group** and select the security group titled "ssh". Do not try to create a new one (you will not be able to). We have forwarded a few ports (they are: 60000-60100, 1024-1048, and 22) that you can use to connect to Jupyter, VScode, Tensorboard, etc. Make sure that the option "**Auto-assign public IP**" to set to enable as in some cases it might be set to disable by default (you might need to click the **edit** button on right top to see the option).
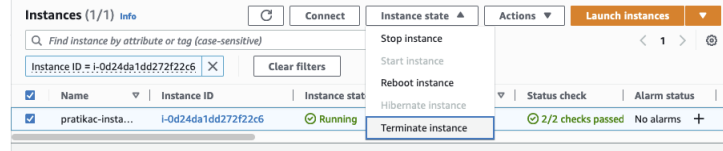


## 4.8 Launch the Instance

Click the **Launch Instance** button and you have a launched GPU instance to play with!

## 4.9 Stop and Start Instance

We have created a script that shuts down idle instances after 1 hour. If you see that your instance has "stopped" you can just start it using the button at the top. **Note: Always stop your instance after use.** Although the provided
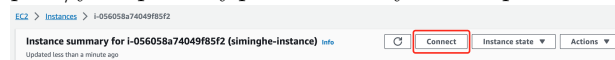
credits should suffice for this course, it is important to use resources carefully. An unstopped instance can quickly deplete the credits. **Stopping** the instance will not delete your data. **Terminating** the instance will delete all your data.
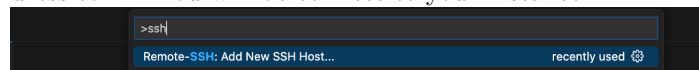


# 5 Using Your Instance

## 5.1 Use Your Instance in VS Code

Click **Connect** and you will be directed to a page called **Connect to instance**. Choose the **SSH Client** tab. You will see an example: "ssh -i 'your-pennkey.pem' address.com". You can use this to access your instance from your laptop. **Important**: you need to change 'your-pennkey.pem' to '/absolute-path/your-pennkey.pem' so that your computer can find the file.



Go to vs code, press **command + shift + P**, type in ssh, and select **add new SSH Host**. Then type in "ssh -i '/absolute-path/your-pennkey.pem' address.com". You will be connect to your instance.



## 5.2 Check the deep learning environment

You can check your environment by opening an terminal in the vs code window and type

```
conda info --env
```

You can activate the environment you want, e.g., pytorch, by

```
conda activate pytorch
```

and run your code by

```
python file-name.py
```

## 5.3 SCP: Secure Copy for File Transfering

You can use the following command to copy file from your computer to your AWS instance and vice versa. You can search for more details online.

```
scp -i '/absolute-path/penn-key.pem' source_path_of_file destination_path_of_file
```

## 5.4   tmux

Training may take a long time. If you accidentally closed the terminal or vs code, the terminal process will be closed which means your training will be killed and lost. You don't want this to happen! If you use tools like tmux, even if you close your terminal, the process will still be ran in the background and you can access it when you come back.Basic commands:

```
# session management
tmux ls (or tmux list-sessions)
tmux new -s session-name
Ctrl-b d Detach from session
tmux attach -t [session name]
tmux kill-session -t session-name

Ctrl-b c Create new window
Ctrl-b d Detach current client
Ctrl-b l Move to previously selected window
Ctrl-b n Move to the next window
Ctrl-b p Move to the previous window
Ctrl-b & Kill the current window
Ctrl-b , Rename the current window
Ctrl-b q Show pane numbers (used to switch between panes)
Ctrl-b o Switch to the next pane
Ctrl-b ? List all keybindings

# moving between windows
Ctrl-b n (Move to the next window)
Ctrl-b p (Move to the previous window)
Ctrl-b l (Move to the previously selected window)
Ctrl-b w (List all windows / window numbers)
Ctrl-b window number (Move to the specified window number, the
default bindings are from 0 -- 9)

# Tiling commands
Ctrl-b % (Split the window vertically)
CTRL-b " (Split window horizontally)
Ctrl-b o (Goto next pane)
Ctrl-b q (Show pane numbers, when the numbers show up type the key to go to that pane)
Ctrl-b { (Move the current pane left)
Ctrl-b } (Move the current pane right)

# Make a pane its own window
Ctrl-b : "break-pane"

# add to ~/.tmux.conf
bind | split-window -h
bind - split-window -v
```

## 5.5   Use Your Instance in Jupyter Notebook

Follow this tutorial! Link