# Machine Learning for Signal Processing
# Independent Component Analysis

Instructor: Bhiksha Raj

# A note on bits..
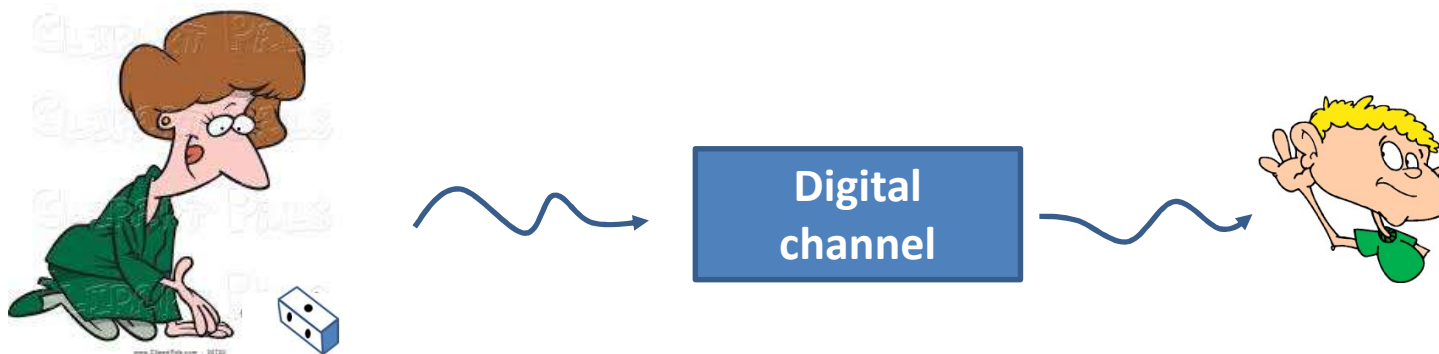
- You flip a coin.  You must inform your friend in the next room about whether the outcome was heads or tails



Digital channel

- How many bits will you have to send?

# A note on bits..

- You roll a four-side dice.  You must inform your friend in the next room about the outcome

**Digital channel**

- How many bits will you have to send?
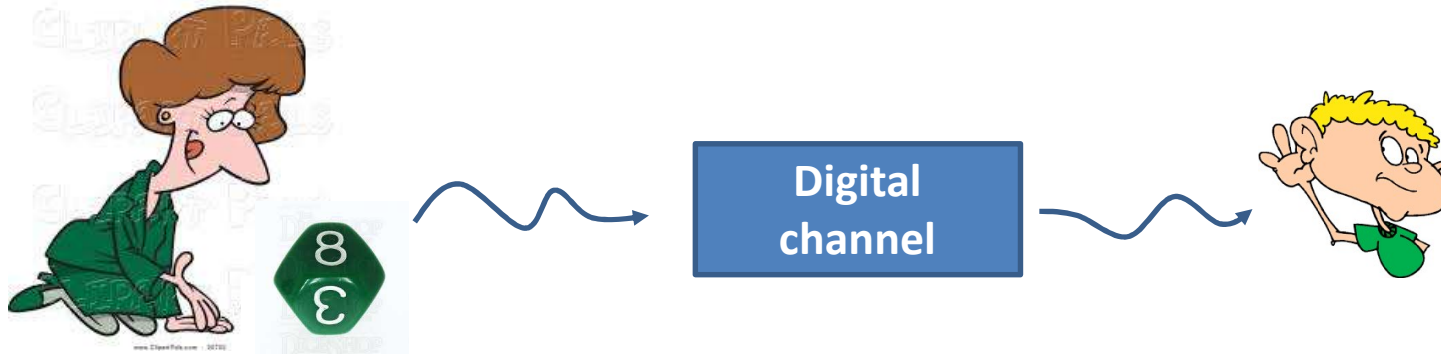
# A note on bits..

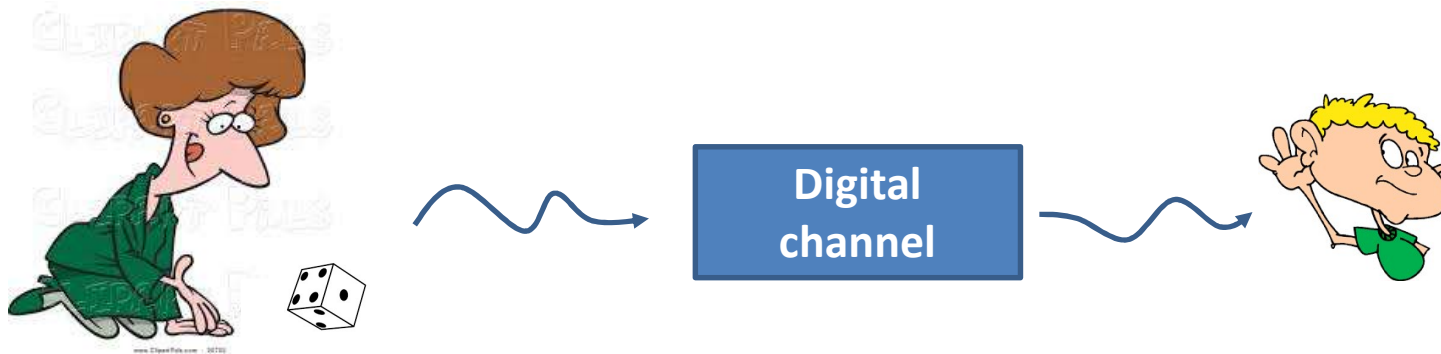- You roll an *eight-sided octahedral* dice.  You must inform your friend in the next room about the outcome



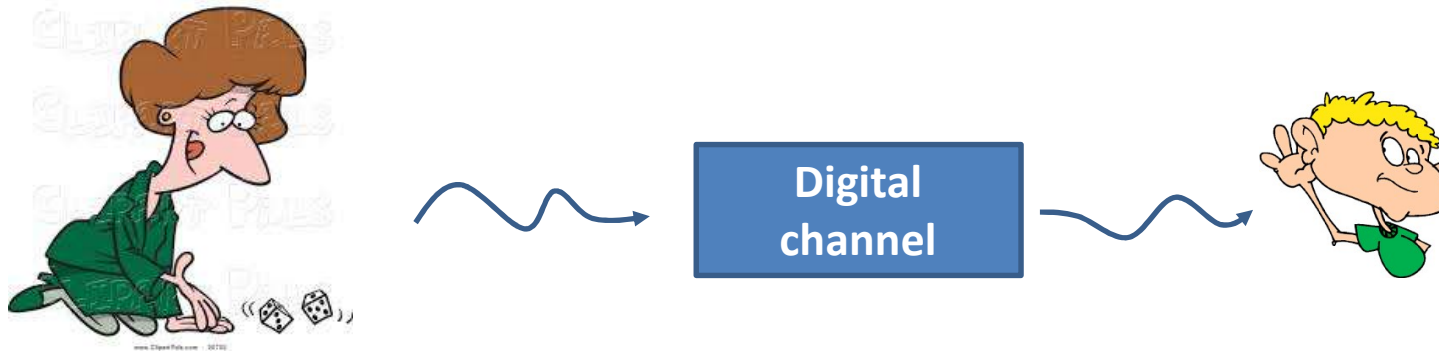- How many bits will you have to send?

# A note on bits..

- You roll a *six-sided* dice.  You must inform your friend in the next room about the outcome

Digital channel

- How many bits will you have to send?
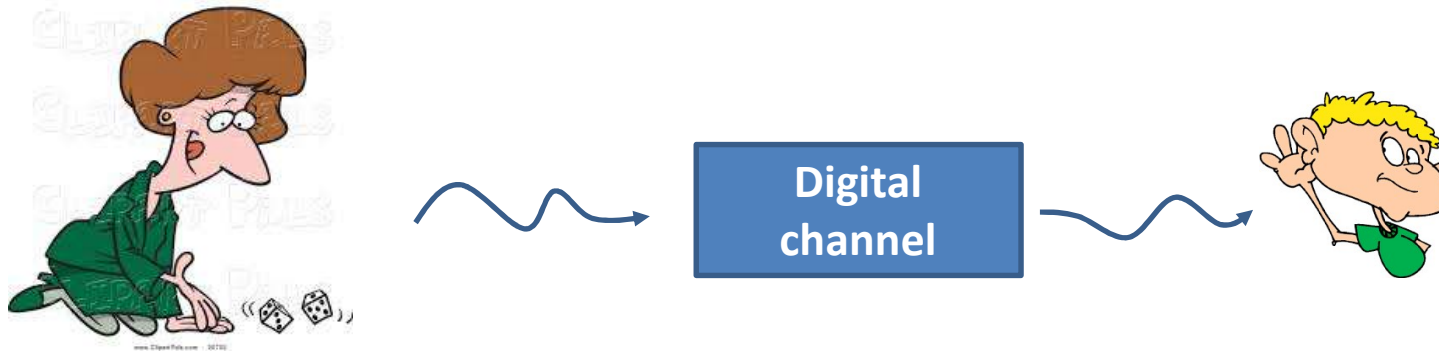
# Batching up 6-sided dice rolls



- Instead of sending individual rolls, you roll the dice *twice*
  - And send the *pair* to your friend
- How many bits do you send *per roll?*

| Roll 1 | Roll 2 |
|--------|--------|
| 1 | 1 |
| 1 | 2 |
| 1 | 3 |
| .. | .. |
| 2 | 1 |
| 2 | 2 |
| .. | .. |
| 6 | 6 |

# Batching up 6-sided dice rolls

**Digital channel**

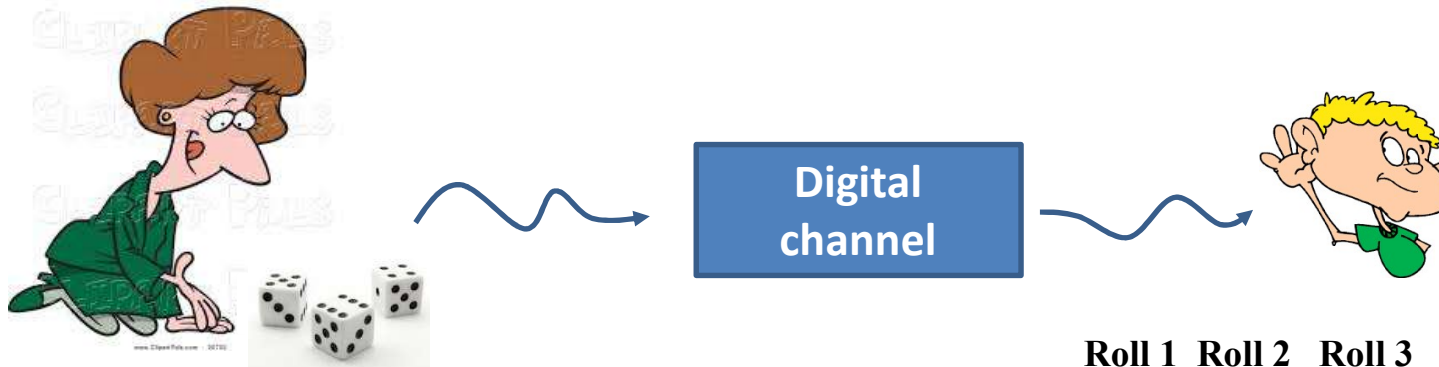| Roll 1 | Roll 2 |
|--------|--------|
| 1 | 1 |
| 1 | 2 |
| 1 | 3 |
| .. | .. |
| 2 | 1 |
| 2 | 2 |
| .. | .. |
| 6 | 6 |

- Instead of sending individual rolls, you roll the dice *twice*
  - And send the *pair* to your friend
- How many bits do you send *per roll?*
- 36 combinations: 6 bits per pair of numbers
  - Still 3 bits per roll

# Batching up 6-sided dice rolls



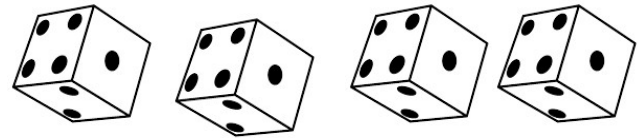| Roll 1 | Roll 2 | Roll 3 |
|--------|--------|--------|
| 1 | 1 | 1 |
| 1 | 1 | 2 |
| .. | .. | .. |
| 1 | 6 | 3 |
| .. |  | .. |
| 2 | 1 | 1 |
| 2 | 1 | 2 |
| .. |  | .. |
| 6 | 6 | 6 |

- Instead of sending individual rolls, you roll the dice **three times**
  - And send the *triple* to your friend
- How many bits do you send *per roll?*
- 216 combinations: 8 bits per triple
  - Still 2.666 bits per roll
  - *Now we're talking!*

# Batching up 6-sided dice rolls

- Batching *four rolls*
    - 1296 combinations
    - 11 bits per outcome (4 rolls)
    - 2.75 bit per roll
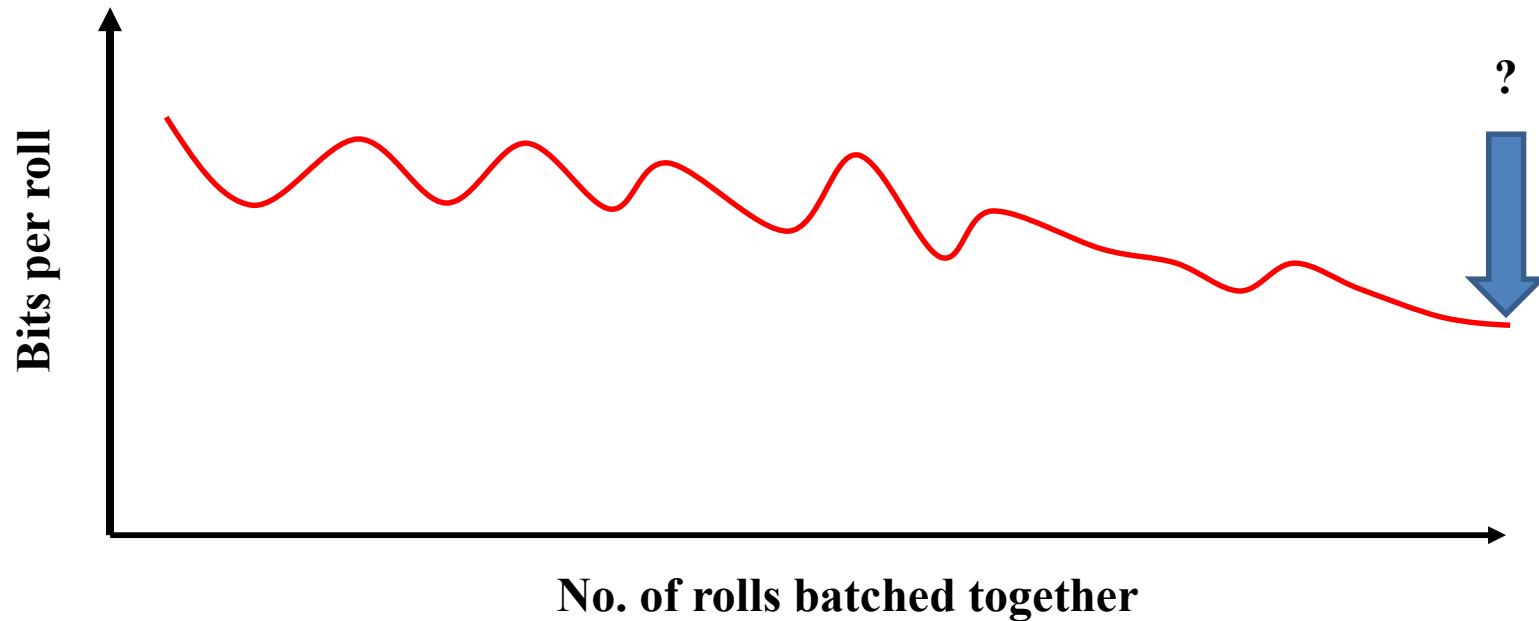
- Batching *five rolls*
    - 7776 combinations
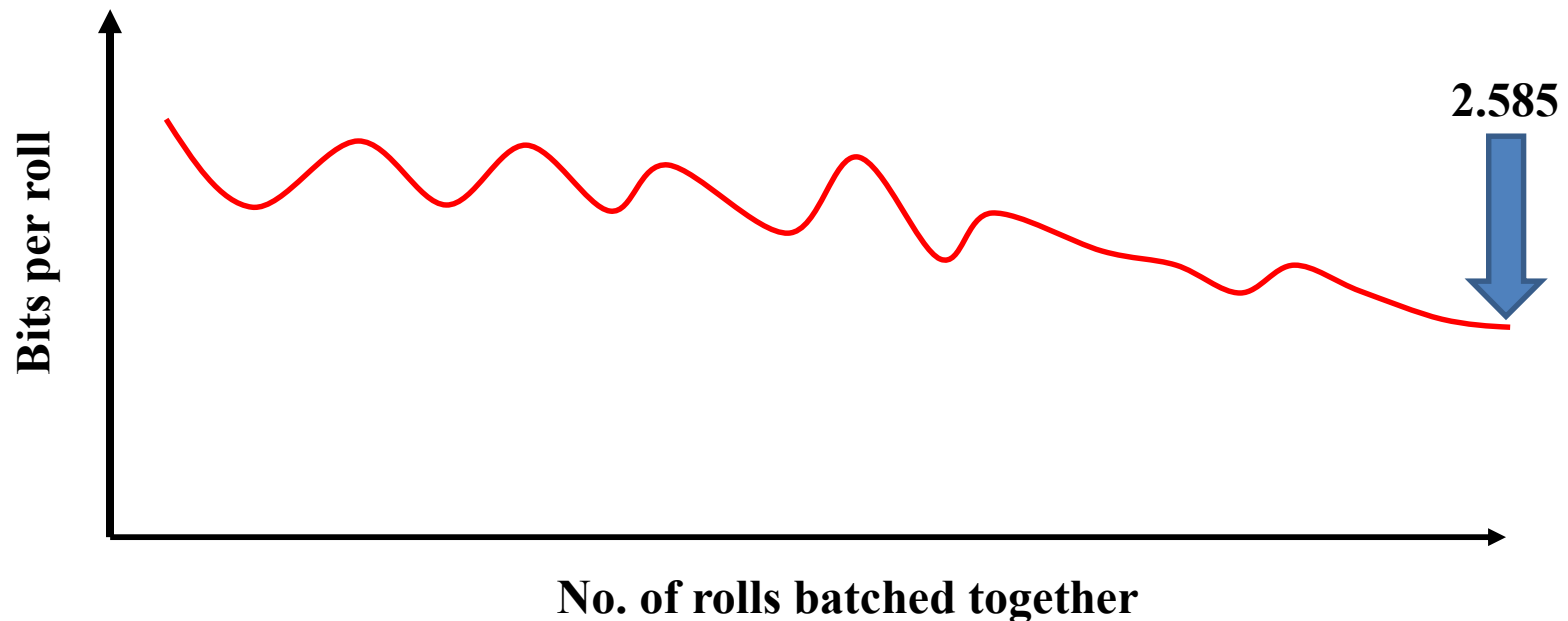    - 13 bits per outcome (5 rolls)
    - 2.6 bits per roll

# Batching up 6-sided dice rolls

Bits per roll

? 

No. of rolls batched together

- Where will it end?

# Batching up 6-sided dice rolls



**2.585**

Bits per roll

No. of rolls batched together

- Where will it end?

- $\lim_{k \to \infty} \frac{[k \log2(6)]}{k} = \log2(6)$ bits per roll in the limit
    - This is the absolute minimum – no simple batching will give you less than these many bits per outcome with this scheme

# Poll 1

# Can we do better?

- A four-sided die needs 2 bits per roll

- But then you find not all sides are equally likely


- P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- *Can you do better than 2 bits per outcome*

# Can we do better?

- You have

P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- You use:

| 1 | 0 |
|---|---|
| 2 | 1 0 |
| 3 | 1 1 0 |
| 4 | 1 1 1 |

  – Note receiver is *never in any doubt as to what they received*
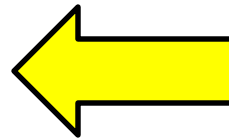
- What is the average number of bits per outcome

# Can we do better?

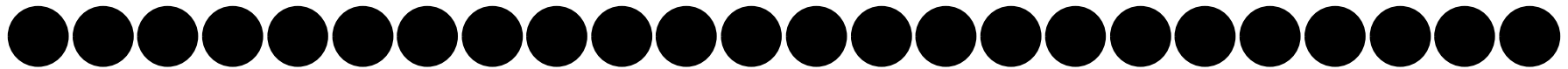- You have

P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- You use:

| 1 | 0 |
|---|---|
| 2 | 1 0 |
| 3 | 1 1 0 |
| 4 | 1 1 1 |

  – Note receiver is *never in any doubt as to what they received*

- How did we know to use three bits here for rows 3 and 4, 2 for row 2 and 1 for row 1?
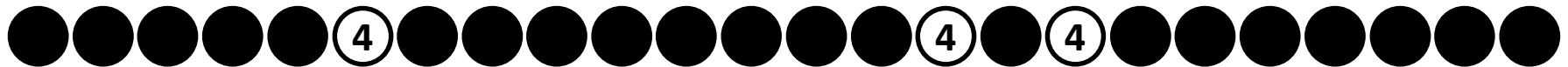
# In a loooong sequence of trials...

●●●●●●●●●●●●●●●●●●●●●●●●

- What fraction of these trials will be "4"?
  - P(4) = 0.125

# In a loooong sequence of trials...

●●●●●④●●●●●●●●④●④●●●●●●

- What fraction of these trials will be "4"?
  - P(4) = 0.125

- From how many alternatives (on average) do we choose the 4
  - From the local perspective of 4

# In a loooong sequence of trials…

●●●●●④●●●●●●●●④●④●●●●●●●

- What fraction of these trials will be "4"?
  - P(4) = 0.125

- From how many alternatives (on average) do we choose the 4
  - From the perspective of 4, you might as well have been rolling an *eight-sided dice*

# In a loooong sequence of trials…

●●●●●●④●●●●●●●●●④●④●●●●●●●●

- What fraction of these trials will be "4"?
  - P(4) = 0.125
- From how many alternatives (on average) do we choose the 4
  - From the perspective of 4 you might as well have been rolling an eight-sided dice
- How many bits to code each instance of 4?
  - When 4 is the outcome of rolls of an 8-sided dice

# In a loooong sequence of trials...

●●●●●④●●●●●●●●●④●④●●●●●●●

- What fraction of these trials will be "4"?
  - P(4) = 0.125
- From how many alternatives (on average) do we choose the 4
  - From the perspective of 4 you might as well have been rolling an eight-sided dice
- How many bits to code each instance of 4?
  - When 4 is the outcome of rolls of an 8-sided dice
- What is the average (expected) number of bits to transmit all instances of 4 in N rolls of the dice?

# In a loooong sequence of trials…



- What fraction of these trials will be "4"?
  - P(4) = 0.125
- From how many alternatives (on average) do we choose the 4
  - From the perspective of 4 you might as well have been rolling an eight-sided dice
- How many bits to code each instance of 4?
  - When 4 is the outcome of rolls of an 8-sided dice
- What is the average (expected) number of bits to transmit all instances of 4 in N rolls of the dice?
  - Average per roll?

# In a loooong sequence of trials...

●●●●●●●●●●●●●●●●●●●●●●●●●●

- What fraction of these trials will be "1"?
  - P(1) = 0. 5

# In a loooong sequence of trials…

①●①①●●①●●●①①●①●①●①①①●●●①

- What fraction of these trials will be "1"?
  - P(1) = 0.5

- From how many alternatives (on average) do we choose the 1
  - From the local perspective of 1

# In a loooong sequence of trials...



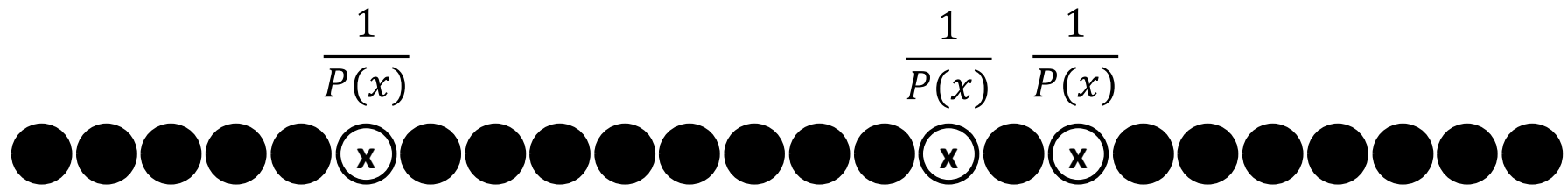- What fraction of these trials will be "1"?

  - P(1) = 0.5

- From how many alternatives (on average) do we choose the 1

  - From the perspective of 1, you might as well have been flipping a coin

# In a loooong sequence of trials…

①●①①●●①●●●①①●①●①●①①①●●●①

- ## What fraction of these trials will be "1"?
  - P(1) = 0.5

- ## From how many alternatives (on average) do we choose the 1
  - From the perspective of 1 you might as well have been flipping a coin

- ## How many bits to code each instance of 1?
  - When 1 is the outcome of a coin toss

# In a loooong sequence of trials…

**①**⚫**①①**⚫⚫**①**⚫⚫⚫**①①**⚫**①**⚫**①**⚫**①①①**⚫⚫⚫**①**

- What fraction of these trials will be "1"?
  - P(1) = 0.5
- From how many alternatives (on average) do we choose the 1
  - From the perspective of 1 you might as well have been flipping a coin
- How many bits to code each instance of 1?
  - When 1 is the outcome of rolls of a coin toss
- What is the average (expected) number of bits to transmit all instances of 1 in N rolls of the dice?

# In a loooong sequence of trials…

①⚫①①⚫⚫①⚫⚫⚫①①⚫①⚫①⚫①①①⚫⚫⚫①

- What fraction of these trials will be "1"?
  - P(1) = 0.5
- From how many alternatives (on average) do we choose the 1
  - From the perspective of 1 you might as well have been flipping a coin
- How many bits to code each instance of 1?
  - When 1 is the outcome of rolls of a coin toss
- What is the average (expected) number of bits to transmit all instances of 1 in N rolls of the dice?
  - Average per roll?

# In a loooong sequence of trials...

- An outcome x has probability P(x)
- From the perspective of x, how many-sided dice is it an outcome of?

# In a loooong sequence of trials...

$$\frac{1}{P(x)} \qquad\qquad \frac{1}{P(x)} \quad \frac{1}{P(x)}$$

- An outcome x has probability P(x)

- From the perspective of x, how many-sided dice is it an outcome of?

- How many bits to code an instance of x?

# In a loooong sequence of trials…

$$\frac{1}{P(x)} \qquad\qquad \frac{1}{P(x)} \quad \frac{1}{P(x)}$$



- An outcome x has probability P(x)

- From the perspective of x, how many-sided dice is it an outcome of?

- How many bits to code an instance of x?

- What is the average (expected) number of bits to transmit instances of x in N rolls of the dice?

# In a loooong sequence of trials…

$$\frac{1}{P(x)} \qquad\qquad \frac{1}{P(x)} \quad \frac{1}{P(x)}$$

⚫⚫⚫⚫⚫ⓧ⚫⚫⚫⚫⚫⚫⚫⚫ⓧ⚫ⓧ⚫⚫⚫⚫⚫⚫⚫

- An outcome x has probability P(x)
- From the perspective of x, how many-sided dice is it an outcome of?
- How many bits to code an instance of x?
- What is the average (expected) number of bits to transmit instances of x in N rolls of the dice?
- Expected number of bits per outcome for *any* outcome?

# In a loooong sequence of trials…

$$\frac{1}{P(x)} \qquad\qquad\qquad \frac{1}{P(x)} \quad \frac{1}{P(x)}$$

●●●●●ⓧ●●●●●●●●●ⓧ●ⓧ●●●●●●●●

- An outcome x has probability P(x)
- From the perspective of x, how many-sided dice is it an outcome of?
- How many bits to code an instance of x?
- What is the average (expected) number of bits to transmit instances of x in N rolls of the dice?
- Expected number of bits per outcome for *any* outcome?
- Average per trial?

# How we do better...

- You have

P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

| 1 | 0 |
|---|---|
| 2 | 1 0 |
| 3 | 1 1 0 |
| 4 | 1 1 1 |

- You use:

  – Note receiver is *never in any doubt as to what they received*

- An outcome with probability $p$ is equivalent to obtaining one of $1/p$ equally likely choices

  – Requires $\log 2 \left( \frac{1}{p} \right)$ bits on average

33

# Entropy



- The average number of bits per symbol required to communicate a random variable over a digitial channel *using an optimal code* is

$$H(p) = \sum_i p_i \log \frac{1}{p_i} = - \sum_i p_i \log p_i$$

- You can't do better
  - Any other code will require more bits
- This is the *entropy of the random variable*

# Poll 2

# A brief review of basic info. theory

T(all),  M(ed), S(hort)…

$$H(X) = \sum_X P(X)[-\log P(X)]$$

- Entropy:  The *minimum average* number of bits to transmit to convey a symbol

**X**

T,  M,  S…

M F  F M..

**Y**

$$H(X,Y) = \sum_{X,Y} P(X,Y)[-\log P(X,Y)]$$

- Joint entropy:  The *minimum average* number of bits to convey sets (pairs here) of symbols

# A brief review of basic info. theory



$$H(X|Y) = \sum_Y P(Y) \sum_X P(X|Y)[-\log P(X|Y)] = \sum_{X,Y} P(X,Y)[-\log P(X|Y)]$$

- Conditional Entropy:  The *minimum average* number of bits to transmit to convey a symbol X, after symbol Y has already been conveyed
  - Averaged over all values of X and Y

And now for something completely different...

# The statistical concept of correlatedness

- Two variables $X$ and $Y$ are correlated if If knowing $X$ gives you an *expected* value of $Y$

- $X$ and $Y$ are uncorrelated if knowing $X$ tells you nothing about the *expected* value of $Y$
  - Although it could give you other information
  - How?

# Correlation vs. Causation

- The consumption of burgers has gone up steadily in the past decade

- In the same period, the penguin population of Antarctica has gone down

Correlation, not Causation
(unless McDonalds has a
 top-secret Antarctica division)

# The concept of *correlation*

- Two variables are correlated if knowing the value of one gives you information about the **expected value** of the other

# A brief review of basic probability

- *Uncorrelated:* Two random variables $X$ and $Y$ are uncorrelated iff:
  - The *average* value of the product of the variables equals the product of their individual averages

- Setup: Each draw produces one instance of $X$ and one instance of $Y$
  - I.e one instance of $(X, Y)$

- $E[XY] = E[X]E[Y]$

- The average value of $Y$ is the same regardless of the value of $X$

# Correlated Variables



- Expected value of $Y$ given $X$:
  - Find average of $Y$ values of all samples at (or close) to the given $X$
  - If this is a function of $X$, $X$ and $Y$ are correlated

# Uncorrelatedness



- Knowing $X$ does not tell you what the *average* value of $Y$ is
  - And vice versa

# Uncorrelated Variables



- The average value of $Y$ is the same regardless of the value of $X$ and vice versa
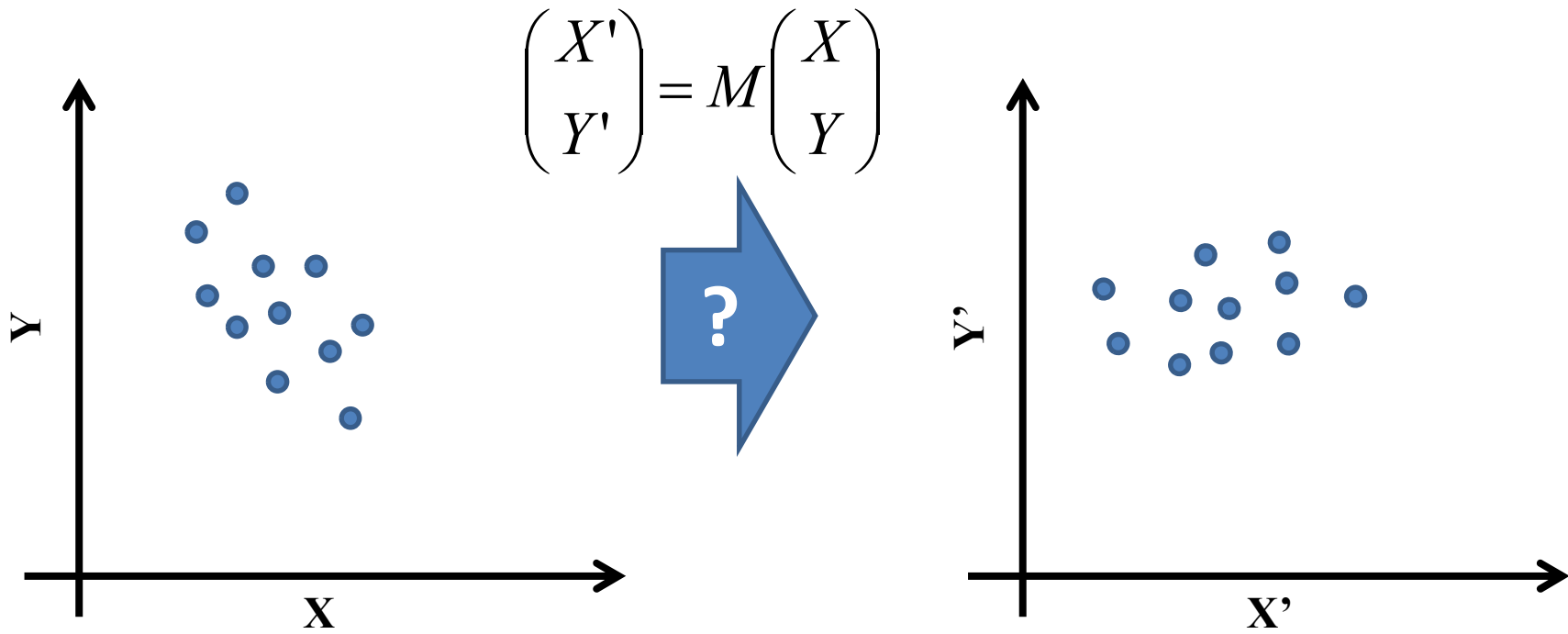
# Uncorrelatedness in Random Variables



- Which of the above represent uncorrelated RVs?

# Benefits of uncorrelatedness..

- Uncorrelatedness of variables is generally considered desirable for modelling and analyses
  - For Euclidean error based regression models and probabilistic models, uncorrelated variables can be separately handled
    - Since the value of one doesn't affect the average value of others
    - Greatly reduces the number of model parameters
  - Otherwise their interactions must be considered

- We will frequently transform correlated variables to make them uncorrelated
  - "Decorrelating" variables

# The notion of *decorrelation*

$$\begin{pmatrix} X' \\ Y' \end{pmatrix} = M \begin{pmatrix} X \\ Y \end{pmatrix}$$



- So how does one transform the correlated variables $(X, Y)$ to the uncorrelated $(X', Y')$

# What does "uncorrelated" mean

**Assuming 0 mean**

**0**

Y'

X'

- $E[X'] = constant$
- $E[Y'] = constant$
- $E[Y'|X'] = constant$
- $E[X'Y'] = E[X']E[Y']$
- All will be 0 for centered data

$$E\left[\begin{pmatrix} X' \\ Y' \end{pmatrix}(X' \quad Y')\right] = E\begin{pmatrix} X'^2 & X'Y' \\ X'Y' & Y'^2 \end{pmatrix} = \begin{pmatrix} E[X'^2] & 0 \\ 0 & E[Y'^2] \end{pmatrix} = diagonal \quad matrix$$

- If $\mathbf{Y}$ is a matrix of vectors, $\mathbf{YY}^{\mathrm{T}}$ = diagonal

# Decorrelation

- Let $\mathbf{X}$ be the matrix of correlated data vectors
  - Each component of $\mathbf{X}$ informs us of the mean trend of other components

- Need a transform $\mathbf{M}$ such that if $\mathbf{Y} = \mathbf{MX}$ such that the covariance of $\mathbf{Y}$ is diagonal
  - $\mathbf{YY}^{\mathrm{T}}$ is the covariance if $\mathbf{Y}$ is zero mean
  - For uncorrelated components, $\mathbf{YY}^{\mathrm{T}} = \mathbf{Diagonal}$
  - $\Rightarrow \mathbf{MXX}^{\mathrm{T}}\mathbf{M}^{\mathrm{T}} = \mathbf{Diagonal}$
  - $\Rightarrow \mathbf{M}.\mathrm{Cov}(\mathbf{X}).\mathbf{M}^{\mathrm{T}} = \mathbf{Diagonal}$

# Decorrelation

- Easy solution:
  - Eigen decomposition of $\mathrm{Cov}(\mathbf{X})$:
  $$\mathrm{Cov}(\mathbf{X}) = \mathbf{E}\Lambda\mathbf{E}^{\mathrm{T}}$$
  - $\mathbf{E}\mathbf{E}^{\mathrm{T}} = \mathrm{I}$
- Let $\mathbf{M} = \mathbf{E}^{\mathrm{T}}$

- $\mathbf{M}\mathrm{Cov}(\mathbf{X})\mathbf{M}^{\mathrm{T}} = \mathbf{E}^{\mathrm{T}}\mathbf{E}\Lambda\mathbf{E}^{\mathrm{T}}\mathbf{E} = \Lambda = \text{diagonal}$

- PCA: $\mathbf{Y} = \mathbf{E}^{\mathrm{T}}\mathbf{X}$

  - Projects the data onto the Eigen vectors of the covariance matrix
  - *Diagonalizes* the covariance matrix
  - "Decorrelates" the data

# PCA

$$\mathbf{X} = w_1 \mathbf{E}_1 + w_2 \mathbf{E}_2$$



- PCA: $\mathbf{Y} = \mathbf{E}^T \mathbf{X}$
  - Projects the data onto the Eigen vectors of the covariance matrix
    - Changes the coordinate system to the Eigen vectors of the covariance matrix
  - *Diagonalizes* the covariance matrix
  - "Decorrelates" the data

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?
- What about if we don't require them to be orthogonal?

# Decorrelating the data



- Are there other decorrelating axes?

- What about if we don't require them to be orthogonal?

- What is special about these axes?

# Poll 3

# The statistical concept of *Independence*

- Two variables X and Y are *dependent* if If knowing X gives you *any information about* Y

- X and Y are *independent* if knowing X tells you nothing at all of Y

# A brief review of basic probability

- ***Independence:*** Two random variables $X$ and $Y$ are independent iff:
  - Their joint probability equals the product of their individual probabilities
- $P(X,Y) = P(X)P(Y)$
- Independence implies uncorrelatedness
  - The average value of $X$ is the same regardless of the value of $Y$
    - $E[X|Y] = E[X]$
  - But uncorrelatedness does not imply independence
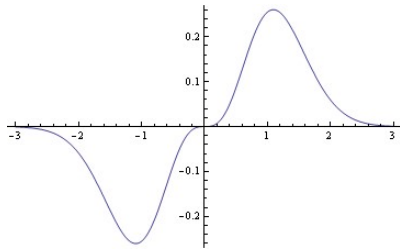
# A brief review of basic probability

- *Independence:* Two random variables $X$ and $Y$ are independent iff:

- The average value of ***any function*** *of $X$* is the same regardless of the value of $Y$
  - Or any function of $Y$

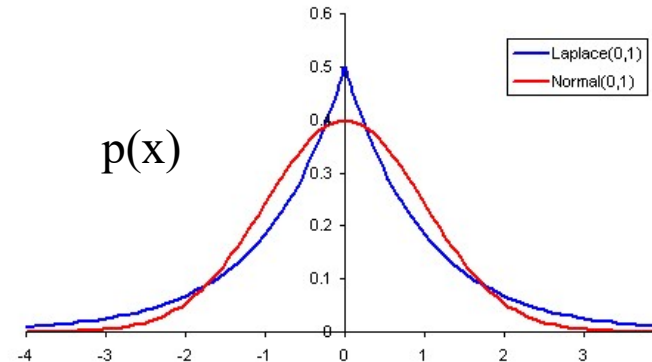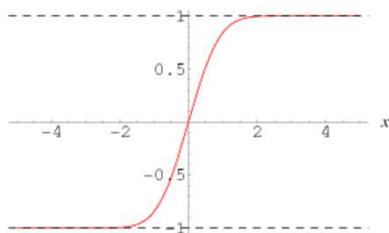- E[f(X)g(Y)] = E[f(X)] E[g(Y)]   for all f(), g()

# Independence



- Which of the above represent independent RVs?
- Which represent uncorrelated RVs?

# A brief review of basic probability

y = f(x)

p(x)

- The expected value of an odd function of an RV is 0 if
  - The RV is 0 mean
  - The PDF is of the RV is symmetric around 0
- **E[f(X)] = 0 if f(X) is odd symmetric**

# A brief review of basic info. theory

- Conditional entropy of $X|Y = H(X)$ if $X$ is independent of $Y$

$$H(X|Y) = \sum_Y P(Y) \sum_X P(X|Y)[-\log P(X|Y)] = \sum_Y P(Y) \sum_X P(X)[-\log P(X)] = H(X)$$

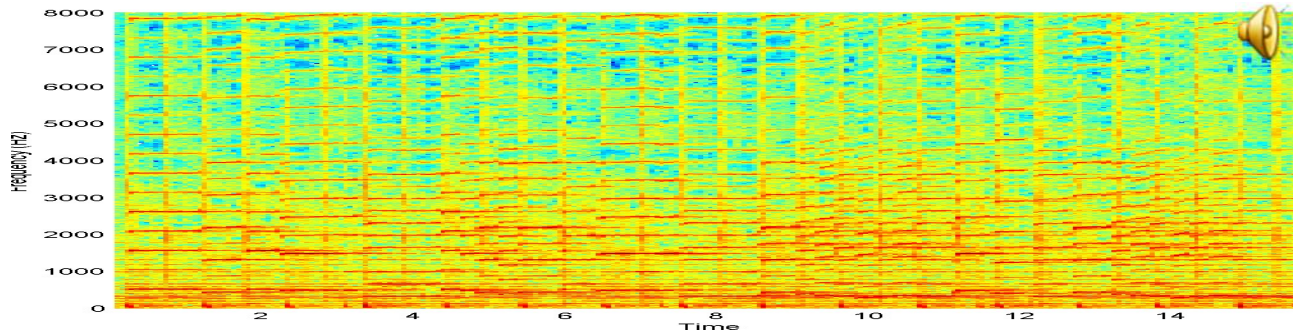- Joint entropy of $X$ and $Y$ is the sum of the entropies of $X$ and $Y$ if they are independent

$$H(X,Y) = \sum_{X,Y} P(X,Y)[-\log P(X,Y)] = \sum_{X,Y} P(X,Y)[-\log P(X)P(Y)]$$

$$= -\sum_{X,Y} P(X,Y)\log P(X) - \sum_{X,Y} P(X,Y)\log P(Y) = H(X) + H(Y)$$

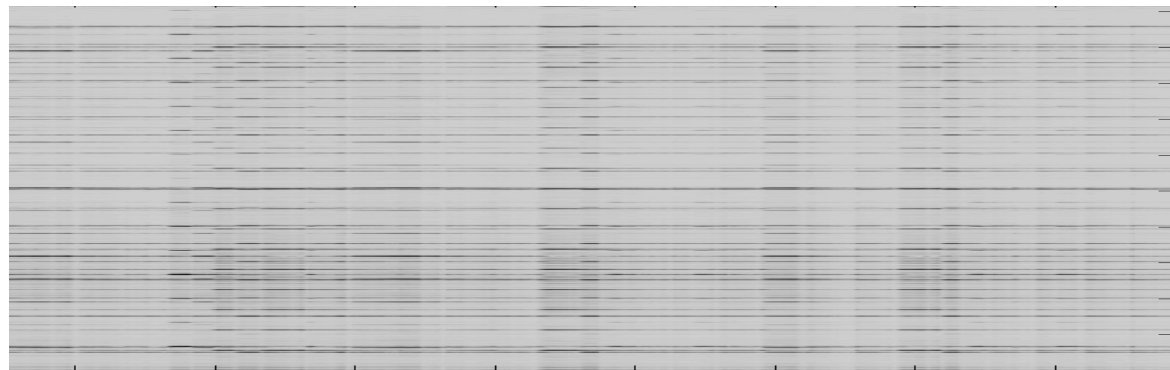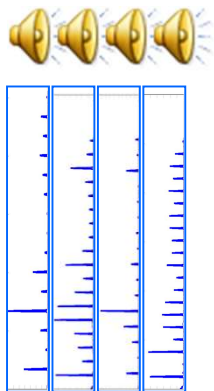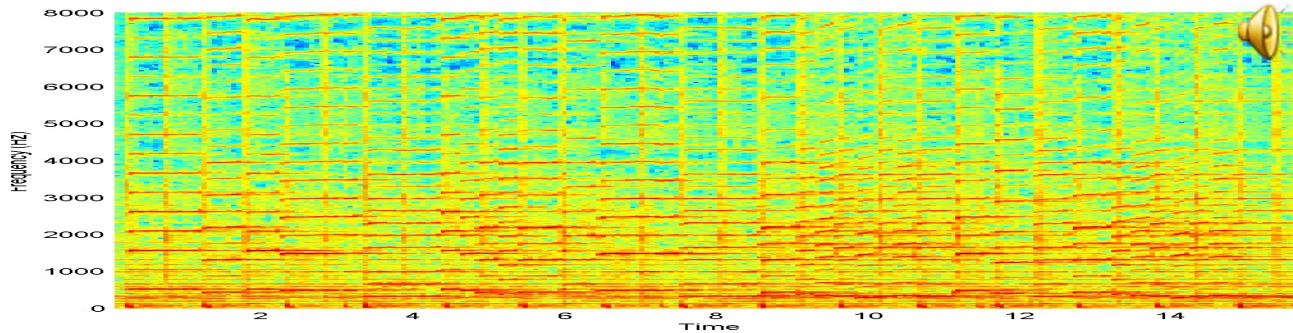# Onward..

# Projection: multiple notes

**M =**



**W =**



- $\mathbf{P} = \mathbf{W} \, (\mathbf{W}^T\mathbf{W})^{-1} \, \mathbf{W}^T$
- Projected Spectrogram = $\mathbf{PM}$

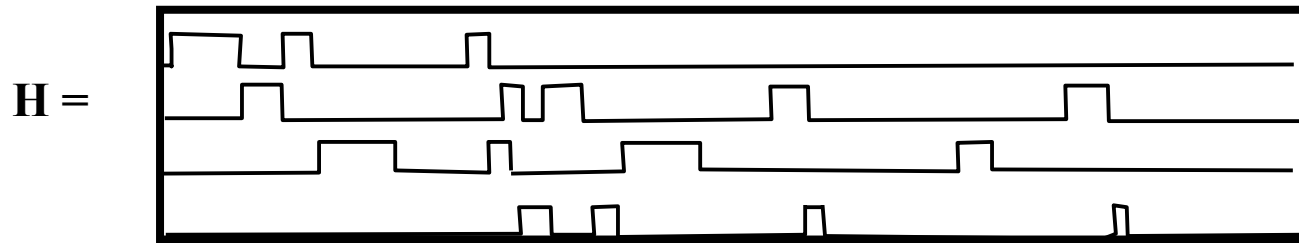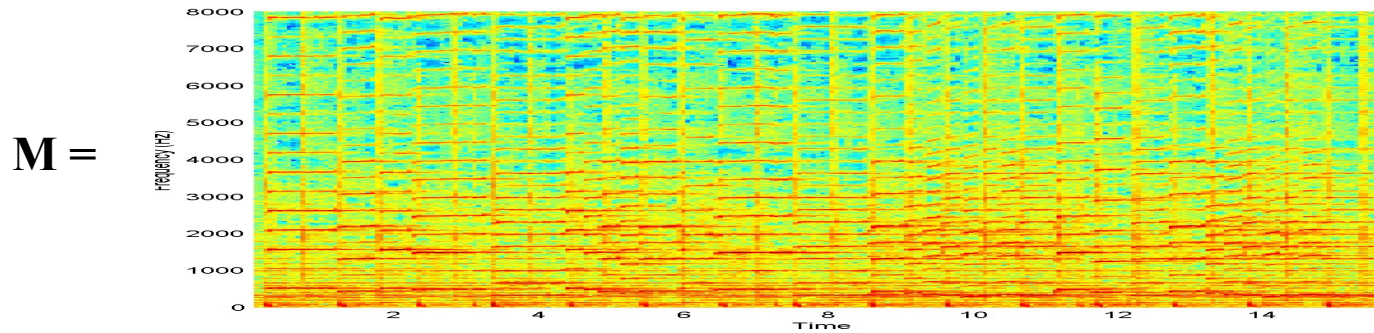# We're actually computing a score

**M =**



$$H = ?$$

**W =**



- **M ~ WH**
- **H = pinv(W)M**

# How about the other way?

**M =**



**H =**



**W =** ?     **U =** ?

- **M ~ WH**     **W = Mpinv(H)**     **U = WH**

# When both parameters are unknown

**H = ?**

**approx(M) = ?**

**W =?**

- Must estimate both $\mathbf{H}$ and $\mathbf{W}$ to best approximate $\mathbf{M}$
- Ideally, must learn *both* the *notes* and *their* transcription!

# A least squares solution

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{W}\mathbf{H}} \|_F^2 + \Lambda(\overline{\mathbf{W}^T}\overline{\mathbf{W}} - \mathbf{I})$$

- Constraint: $\mathbf{W}$ is orthogonal
  - $\mathbf{W}^T\mathbf{W} = \mathbf{I}$
- The solution: $\mathbf{W}$ are the Eigen vectors of $\mathbf{MM}^T$
  - PCA!!

- $\mathbf{M} \sim \mathbf{WH}$ is an approximation
- Also, the rows of $\mathbf{H}$ are *decorrelated*
  - Trivial to prove that $\mathbf{HH}^T$ is diagonal

# PCA

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{W}\mathbf{H}} \|_F^2$$

$$\mathbf{M} \approx \mathbf{W}\mathbf{H}$$

$$\mathbf{W}\mathbf{W}^{\mathbf{T}} = \text{Diagonal} \quad \text{OR} \quad \mathbf{H}\mathbf{H}^{\mathbf{T}} = \text{Diagonal}$$
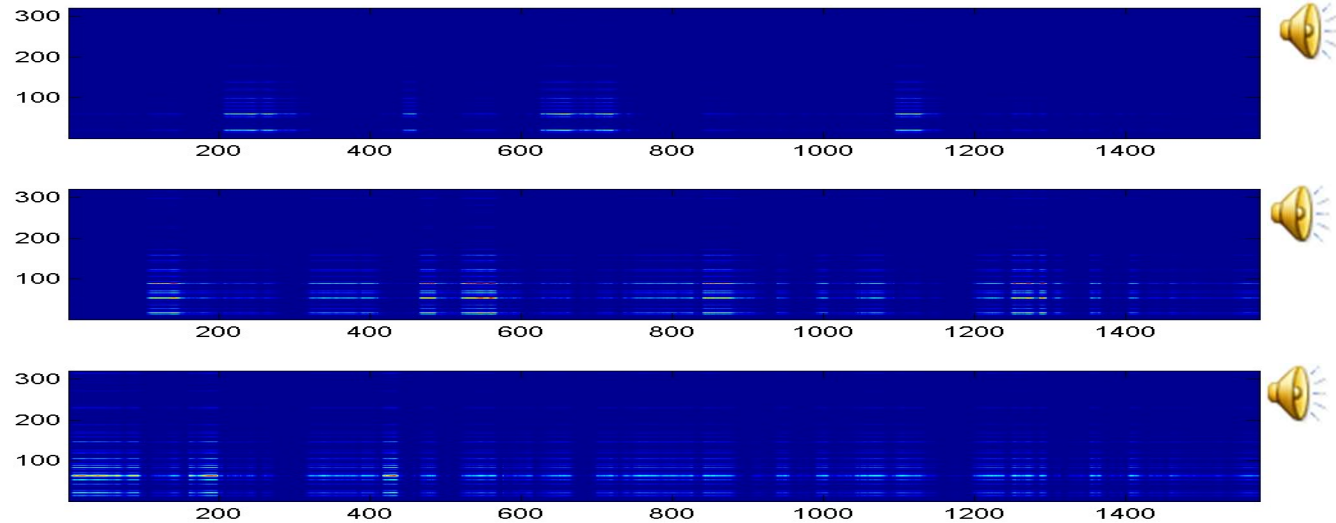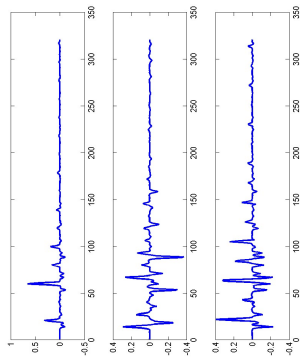
**The conditions are equivalent**

- The columns of **W** are the bases we have learned
  - The linear "building blocks" that compose the music
- They represent "learned" notes
  - $\mathbf{w}_i \mathbf{h}_i$ is the contribution of the ith note to the music
    - $\mathbf{w}_i$ is the ith column of **W**
    - $\mathbf{h}_i$ is the ith row of **H**

# So how does that work?



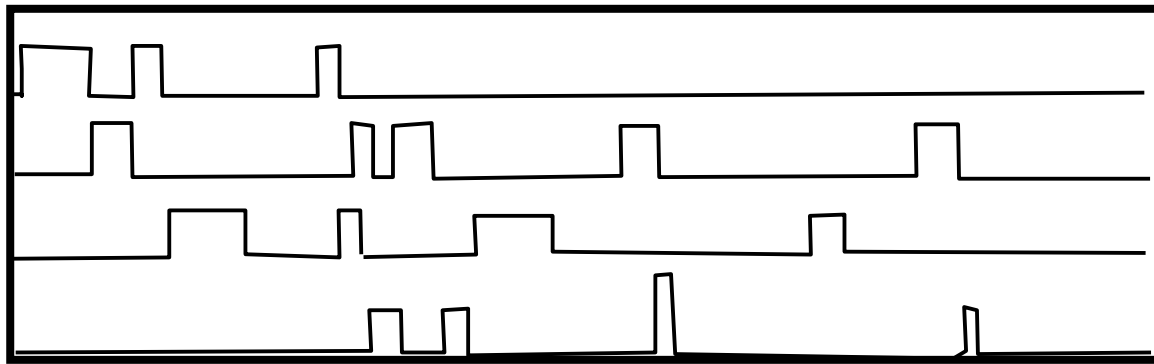- There are 12 notes in the segment, hence we try to estimate 12 notes..

# So how does that work?



- There are 12 notes in the segment, hence we try to estimate 12 notes..
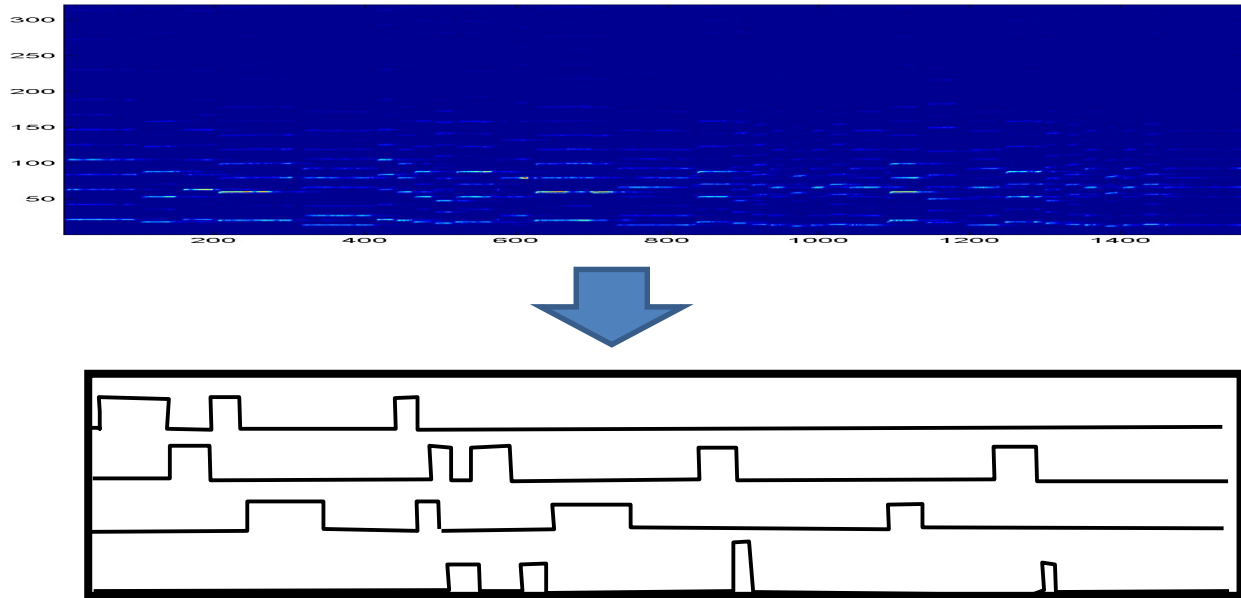- Results are not good

# PCA through decorrelation of notes

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{H}} \|_F^2 + \Lambda(\overline{\mathbf{H}\mathbf{H}}^T - \mathbf{D})$$



- Different constraint: Constraint $\mathbf{H}$ to be decorrelated
  - $\mathbf{H}\mathbf{H}^\mathrm{T} = \mathbf{D}$
- This will result exactly in PCA too
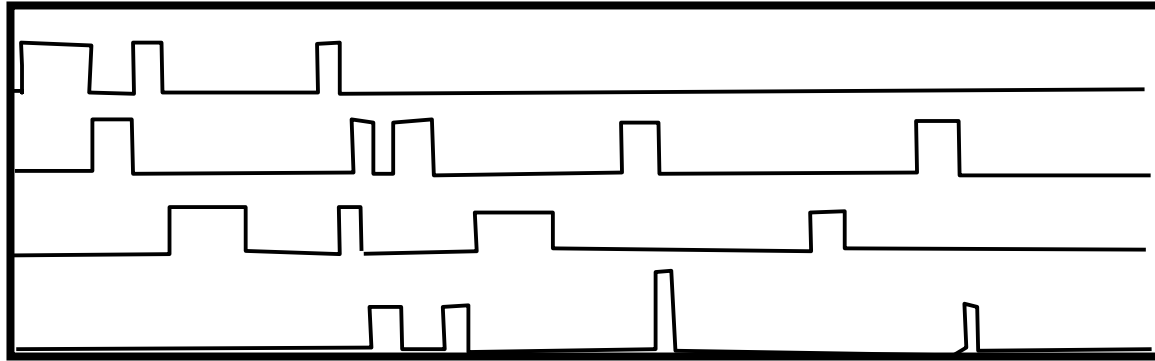- Decorrelation of $\mathbf{H}$ Interpretation: What does this mean?

# Decorrelation



- Alternate view: Find a matrix $\mathbf{B}$ such that the rows of $\mathbf{H}=\mathbf{BM}$ are uncorrelated

- Will find $\mathbf{B} = \mathbf{W}^{\mathrm{T}}$

- $\mathbf{B}$ is the *decorrelating matrix* of $\mathbf{M}$

# Poll 4

# What *else* can we look for?



- Assume: The "transcription" of one note does not depend on what else is playing
  - Or, in a multi-instrument piece, instruments are playing independently of one another
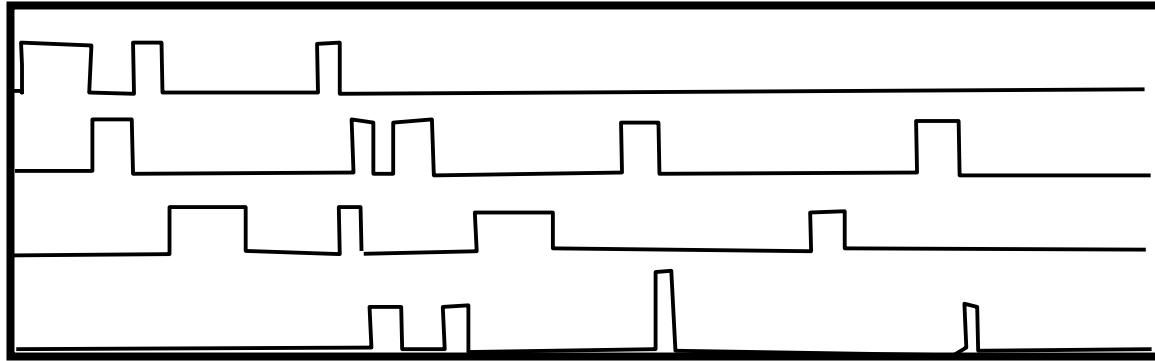- Not strictly true, but still..

# What *else* can we look for?



- Assume: The "transcription" of one note does not depend on what else is playing
  - Or, in a multi-instrument piece, instruments are playing independently of one another

- **Attempting to find statistically independent components of the mixed signal**
  - *Independent Component Analysis*

# Formulating it with Independence

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{W}\mathbf{H}} \|_F^2 + \Lambda(\textit{rows of } \mathbf{H} \textit{ are independent})$$

- Impose statistical independence constraints on decomposition

# Next Class

- Independent Component Analysis

- By Yinghao Ma