# Probability

$P(A) = $ "probability of event $A$"

$P(A|B) = $ "probability of event $A$ <u>given</u> event $B$"

<span style="color:red">↖ "given"</span>

$P(\bar{A}) = $ "probability of event $A$ <u>not</u> occurring"

# Basic Rules:

* sum of all probabilities in an event space $= 1$
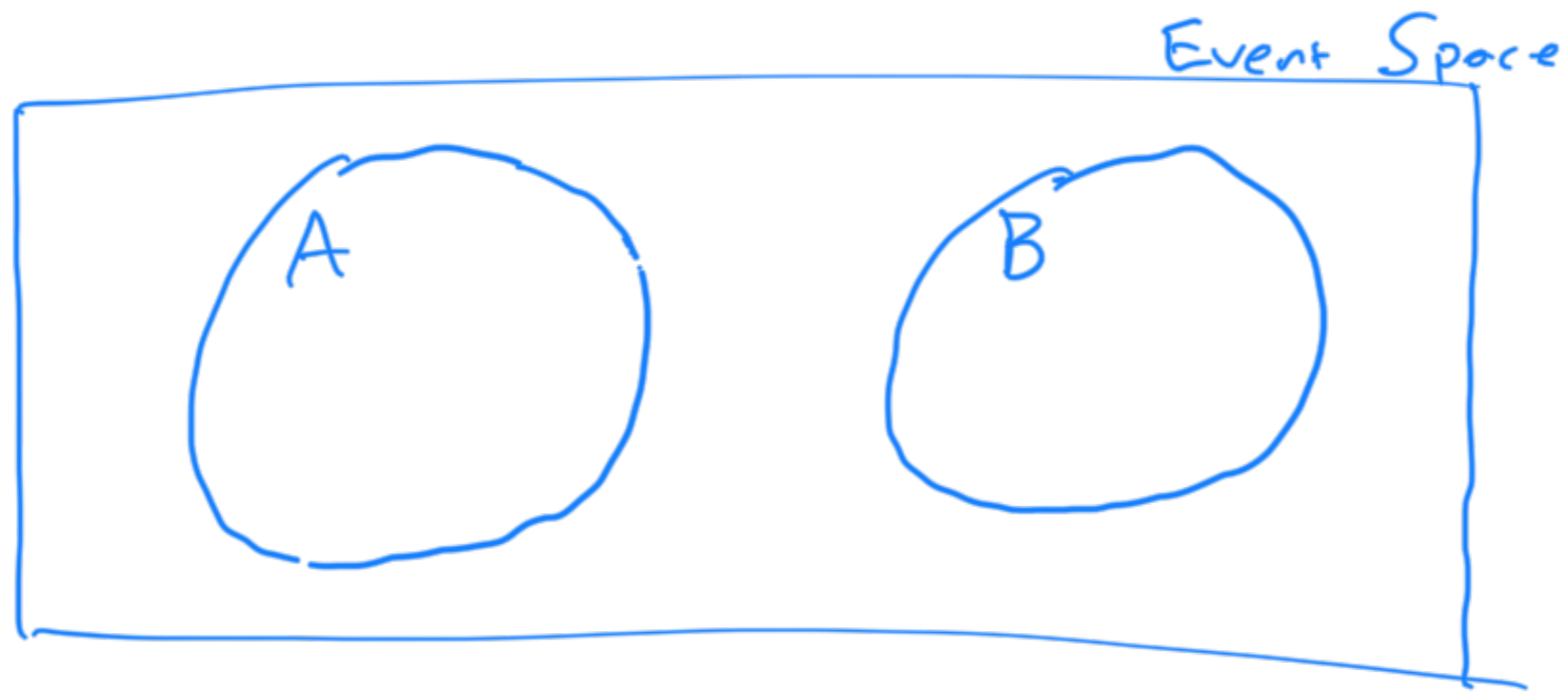
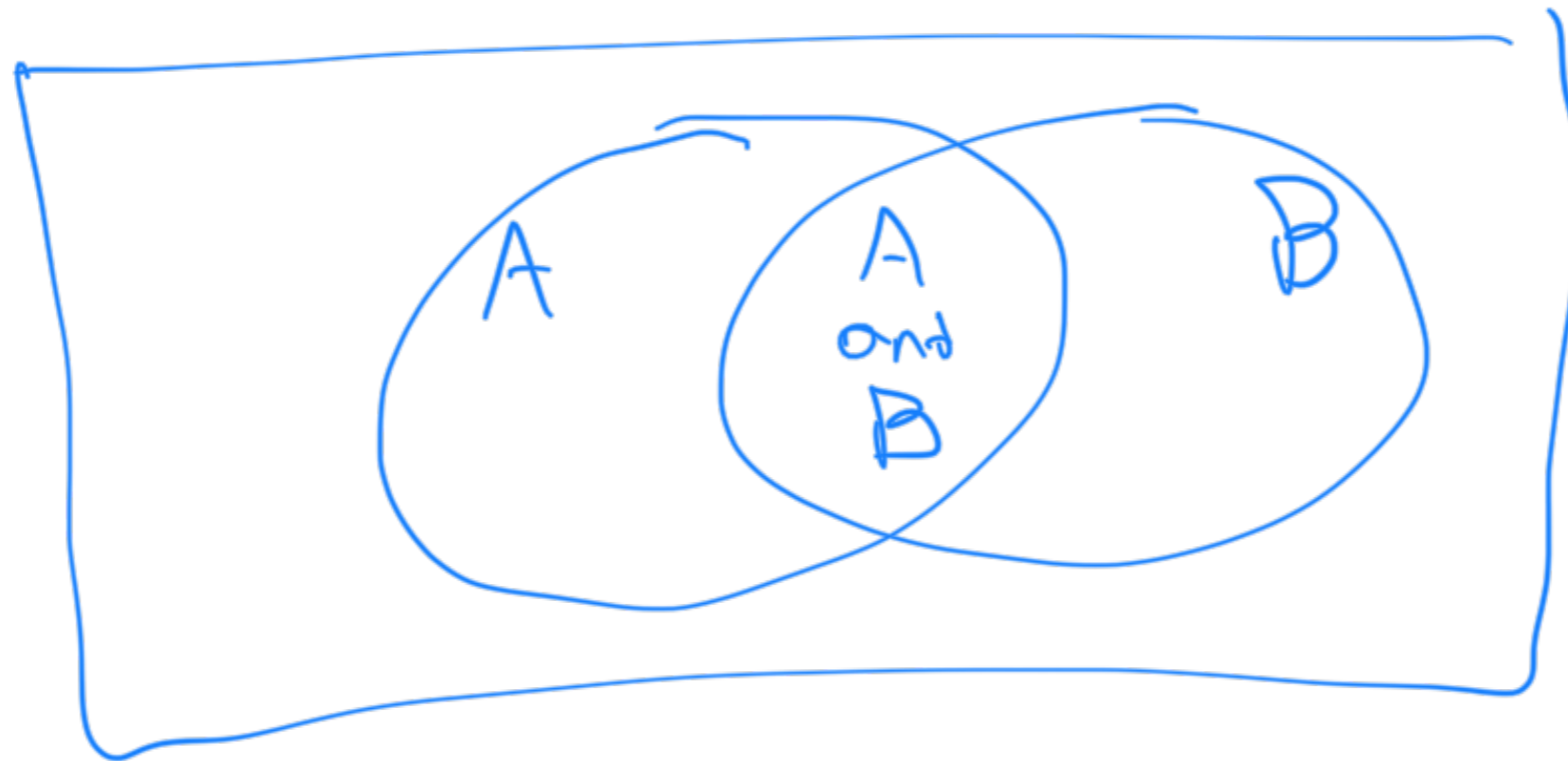* $0 \leq P(x) \leq 1$

* $P(\bar{A}) = 1 - P(A)$

  $P(A) = 1 - P(\bar{A})$

* $P(A \text{ or } B) = P(A) + P(B)$

  <u>iff</u> $A$ and $B$ are <u>mutually</u> <u>exclusive</u>

Event Space

A

B

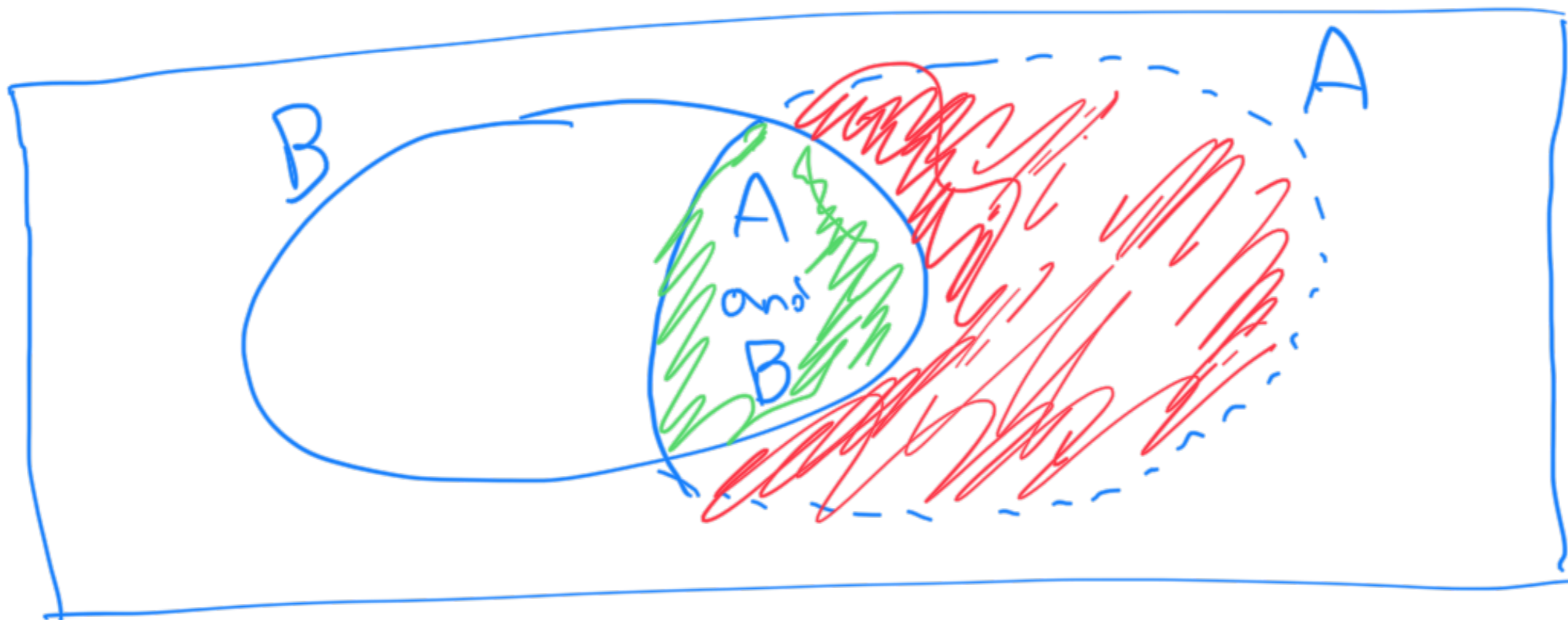$*$ $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

A

A and B

B

$*$ $P(A \text{ and } B) = \dfrac{P(A)P(B|A)}{P(A)P(B)} = P(B)P(A|B)$

$*$ $P(A \text{ and } B) = P(A)P(B)$

iff A and B are independent

* Bayes' Rule:

$$P(A|B) = \frac{P(A \text{ and } B)}{P(B)}$$



= event A that is relevant given B
= event A that is irrelevant given B

# Email Spam Categorization using Probability

In Excel:

| Text | Spam? |
|------|-------|
| Buy this pill to ... | 1 |
| Free lottery Tickets! Just ... | 1 |
| Can I send you $1 million? ... | 1 |
| ... ... | 1 |
| ... ... | 1 |
| ... ... | 1 |
| Dear Peter, Can you help me with my HW? | 0 |
| Check out this concert on ... | 0 |
| ... ... | 0 |
| ... ... | 0 |

# Step 2

Convert text into numerical representation:

   **1** is word appears in text

   0 is word doesn't appear in text

(this table doesn't correspond to step 1's table)

| hello | Vicodin | ⋯ | Spam? |
|:---:|:---:|:---:|:---:|
| 0 | 1 | | 1 |
| 1 | 1 | | 1 |
| 0 | 0 | | 1 |
| 1 | 1 | | 1 |
| 1 | 1 | | 1 |
| 0 | 0 | | 1 |
| | | | 1 |

$$
\begin{array}{c|c|c}
1 & 0 & \dfrac{1}{0} \\
1 & 0 & 0 \\
1 & 1 & 0 \\
1 & 1 & 0 \\
\end{array}
$$

## Step 3

Calculate the relevant probabilities for Bayes' Rule

$$P(A \mid B) = \frac{P(A \text{ and } B)}{P(B)} = \frac{P(B \mid A)\,P(A)}{P(B)}$$

$$P\left(\text{Spam} \mid \text{each word}\right) = \frac{P\left(\text{word 1} \mid \text{spam}\right) \cdot P\left(\text{word 2} \mid \text{Spam}\right) \cdot \cdots \cdot P\left(\text{Spam}\right)}{P\left(\text{seeing all words}\right)}$$

$$P\left(\text{not Spam} \mid \text{each word}\right) = \frac{P\left(\text{word 1} \mid \text{not spam}\right) \cdot P\left(\text{word 2} \mid \text{not spam}\right) \cdot \cdots \cdot P\left(\begin{array}{c}\text{not}\\\text{Spam}\end{array}\right)}{P\left(\text{seeing all words}\right)}$$

$$P(\text{seeing} \quad \text{word})$$

To categorize as spam, $P\left(\text{spam} \mid \text{each word}\right) > P\left(\begin{smallmatrix}\text{not}\\\text{spam}\end{smallmatrix} \mid \text{each word}\right)$

To categorize or $\begin{smallmatrix}\text{not}\\\text{spam}\end{smallmatrix}$, $P\left(\text{spam} \mid \text{each word}\right) \leq P\left(\begin{smallmatrix}\text{not}\\\text{spam}\end{smallmatrix} \mid \text{each word}\right)$

$$P\left(\text{"hello"} \mid \text{spam}\right) = \frac{\text{number of spam emails with "hello"}}{\text{number of spam emails}} = 50\%$$

$$P\left(\text{"Vicodin"} \mid \text{spam}\right) = \frac{4}{6} = 66.7\%$$

$$P\left(\text{spam}\right) = \frac{6}{10} = 60\%$$

$$P\left(\text{"hello"} \mid \begin{smallmatrix}\text{not}\\\text{spam}\end{smallmatrix}\right) = \frac{4}{4} = 100\%$$

$$P\left(\text{"Vicodin"} \mid \begin{smallmatrix}\text{not}\\\text{spam}\end{smallmatrix}\right) = \frac{2}{4} = 50\%$$

$$P\left(\text{not spam}\right) = 40\%$$

# Step 4

Use these probabilities to categorize new, unseen emails.

Example new email: "Hello! Buy my Vicodin."

$$P(Spam \mid words) \propto P(\text{"hello"} \mid Spam) \, P(\text{"Vicodin"} \mid Spam) \, P(Spam) = 10\%$$

$$P(\text{not} \atop Spam \mid words) \propto P(\text{"hello"} \mid \text{not} \atop Spam) \, P(\text{"Vicodin"} \mid \text{not} \atop Spam) \, P(\text{not} \atop Spam)$$

$$= 100\% \cdot 50\% \cdot 40\% = 20\%$$

Since $P(Spam \mid words) > P(\text{not} \atop Spam \mid words)$,
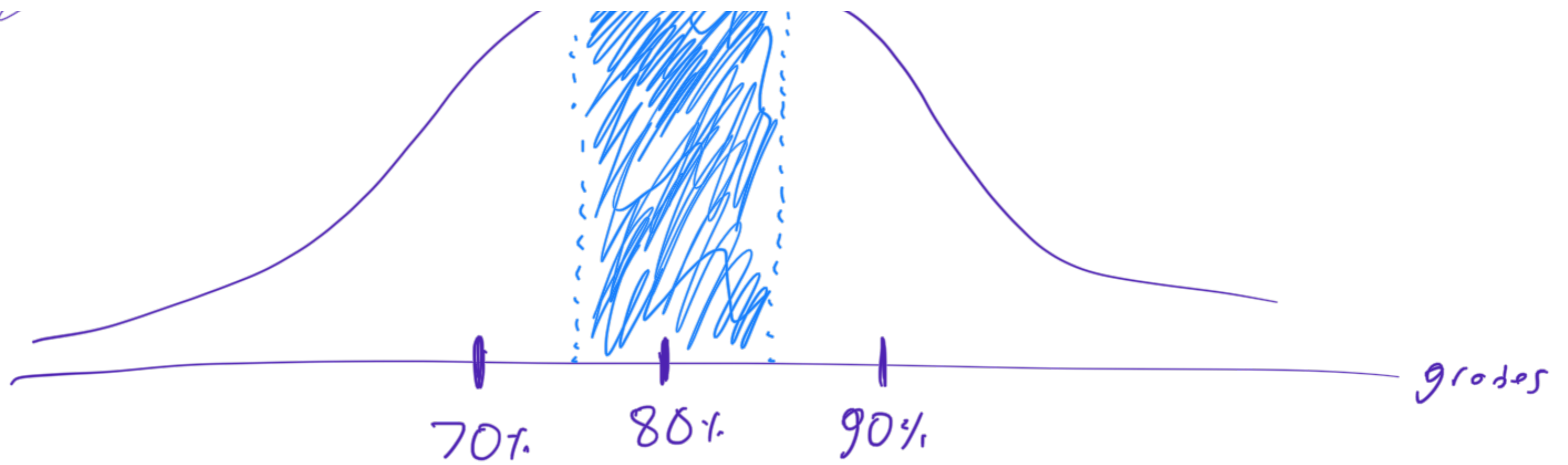
Put email in spam folder.

Example new email: "Hello! Let's go to lunch."

$$P(spam \mid words) \propto P(\text{``hello''} \mid spam) P(\text{not ``vicodin''} \mid spam) P(spam)$$

$$= 50\% \cdot 33.3\% \cdot 60\%$$

$$= 10\%$$

$$P(\text{not spam} \mid words) \propto P\left(\text{``hello''} \mid \begin{array}{c} not \\ spam \end{array}\right) P\left(\begin{array}{c} not \\ \text{``vicodin''} \end{array} \middle| \begin{array}{c} not \\ spam \end{array}\right) \cdot P(\text{not spam})$$

$$= 100\% \cdot 50\% \cdot 40\%$$

$$= 20\%$$

Since $P(\text{not spam} \mid words) > P(spam \mid words)$,

don't put email in spam folder

70%    80%    90%    grades

$$P\left(75\% < \text{grade} < 85\%\right)$$

$$= \int_{75}^{85} \text{Bell Curve } dx$$