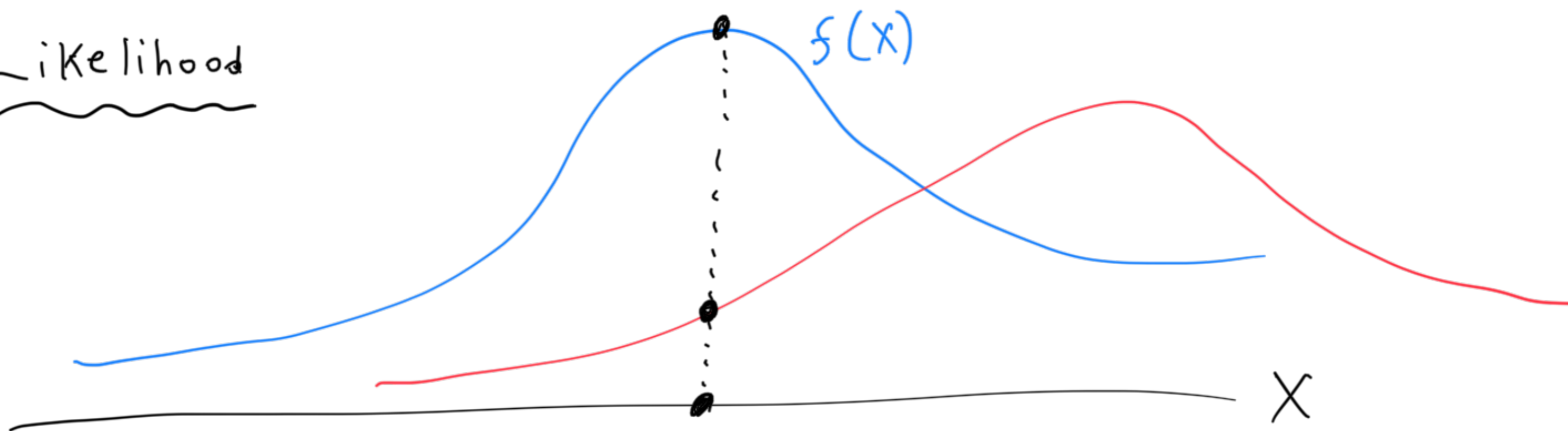


Day 10: Maximum Likelihood Estimation and Non-Parametric Learning Intro

Likelihood



$$L(\theta | X) = f(x_1 | \theta) \times \dots \times f(x_n | \theta)$$

$$\log L(\theta | X) = \log f(x_1 | \theta) + \dots + \log f(x_n | \theta)$$

$\log(f(x))$ increases/decreases same as $f(x)$

Maximum Likelihood Estimation (MLE)

$$\frac{\partial L(\theta | X)}{\partial \theta} = 0$$

Solve for θ

regular
optimization
problem

In practice, people solve:

$$\frac{\partial \log L(\theta | X)}{\partial \theta} = 0 \quad \text{for } \theta$$

Because:

* many fractions multiplied will underflow

$$* \operatorname{argmax}_{\theta} f(x) = \operatorname{argmax}_{\theta} \log f(x)$$

$$\operatorname{argmin}_{\theta} f(x) = \operatorname{argmin}_{\theta} \log \mathcal{L}(x)$$

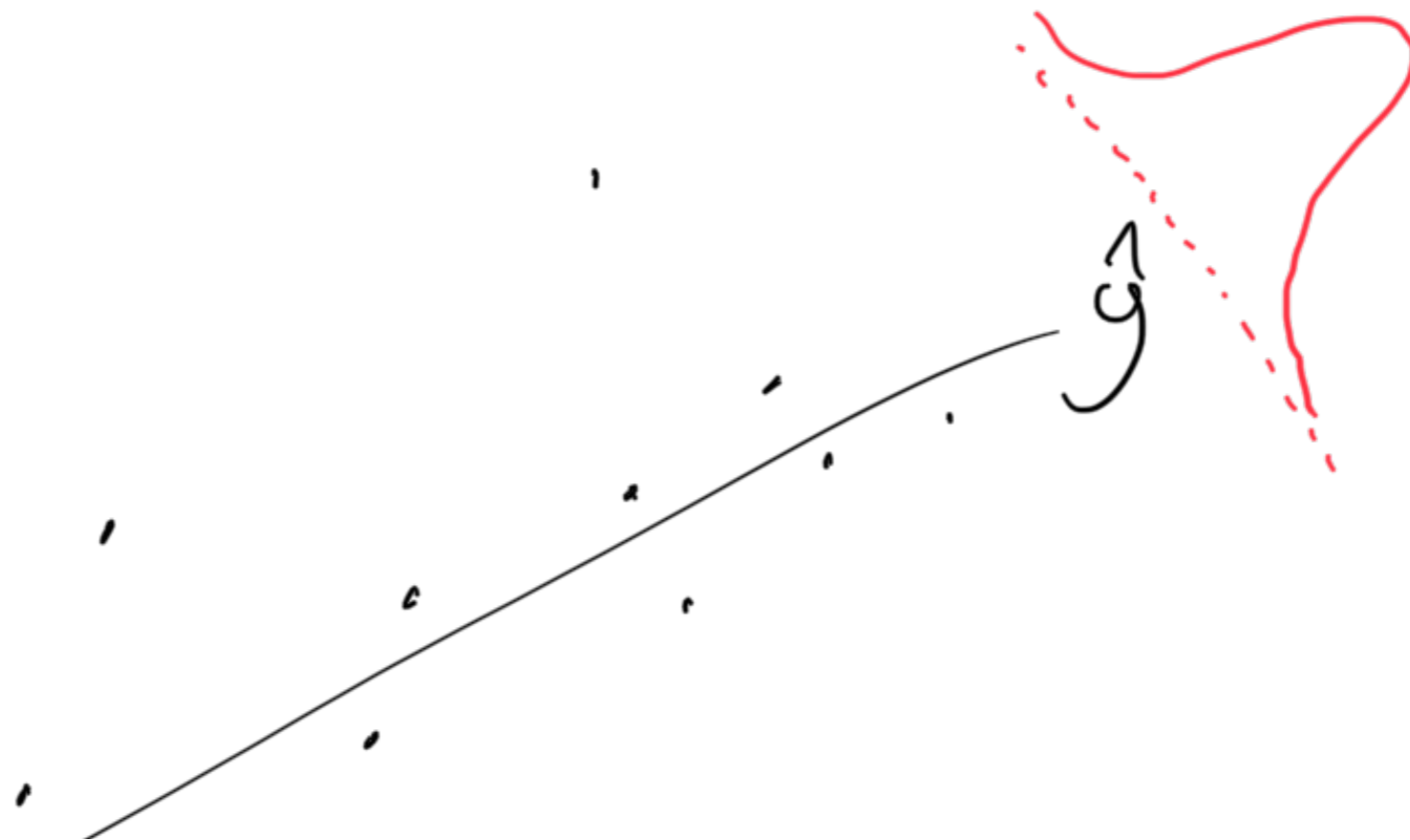
$\left(\operatorname{argmin}_{\theta} f(x) = \text{"value of } \theta \text{ which"} \right.$
 $\left. \text{minimizes } f(x) \right)$

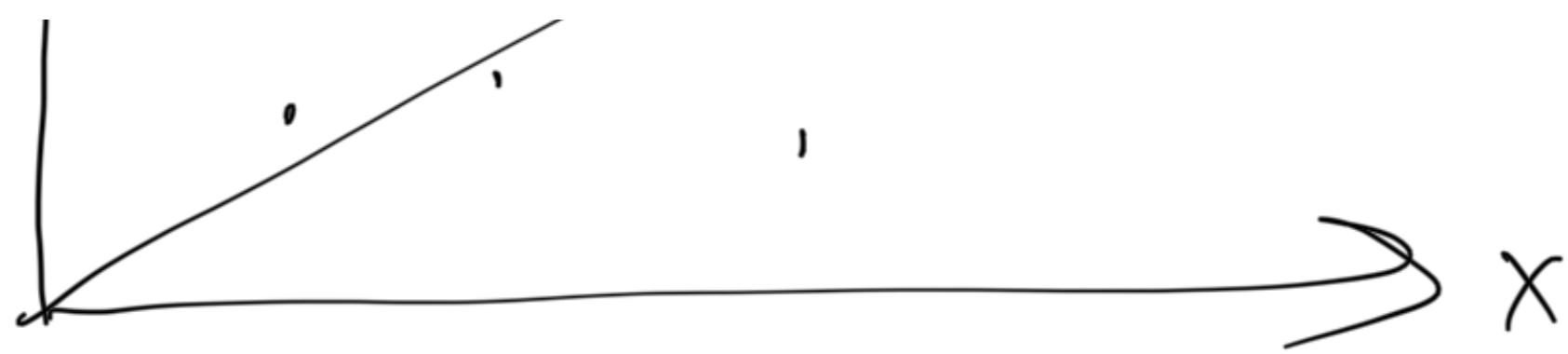
Big Takeaway for Linear Regression:

Maximizing likelihood = minimizing MSE

Proof:

y





$$\hat{y} = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$y = \hat{y} + \epsilon$$

$$\epsilon \sim N(0, \sigma^2)$$

the true data is the predicted line plus some error

assume the error is normally distributed with mean 0 and variance σ^2

So:

$$y \sim N(\hat{y}, \sigma^2)$$

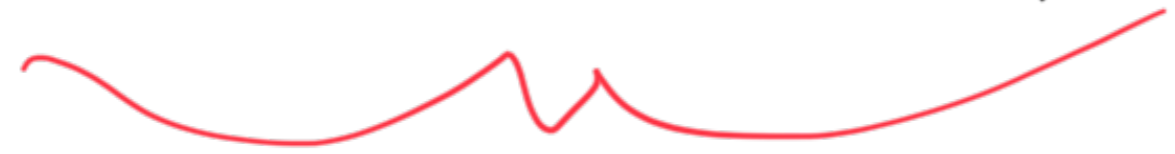
$L(\hat{y}, \sigma^2 | y)$

log-likelihood =

$$\frac{-(y - \hat{y})^2}{2}$$

$$\sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \hat{y}_i)^2}{2\sigma^2}} \right)$$

$$= n \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) - \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{2\sigma^2}$$



Constant in terms of θ



$\text{argmin } f(x) = \text{argmin } C \cdot f(x)$

Solving $\frac{\partial(\log \text{ likelihood})}{\partial \theta} = 0$ for θ

is equivalent to

$$-\sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Solving

$$\frac{\partial \left(- \sum_{i=1}^n \ln(\dots) \right)}{\partial \theta} = 0 \quad \text{for } \theta$$

"Maximum likelihood estimate for θ "

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \left[- \sum_{i=1}^n (y - \hat{y})^2 \right]$$

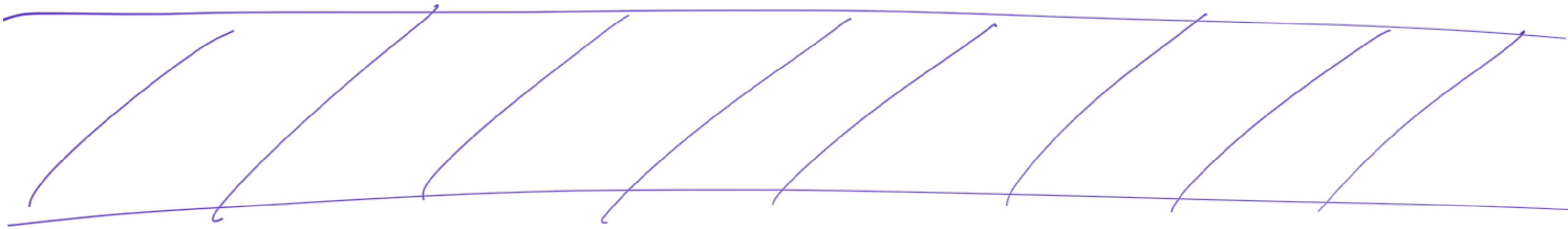
$$= \operatorname{argmin}_{\theta} \left[\sum_{i=1}^n (y - \hat{y})^2 \right]$$

Sum of square errors
(SSE)

$$MSE = \frac{1}{n} SSE$$

Constant

$$\underset{\theta}{\operatorname{argmin}} \text{MSE} = \underset{\theta}{\operatorname{argmin}} \text{SSE}$$



$$L(p) = \prod_{i=1}^n \underbrace{p^{y_i} (1-p)^{(1-y_i)}}_{\text{Bernoulli PMF}}$$

$$\begin{aligned} \log L(p) &= \log \left(\prod_{i=1}^n p^{y_i} (1-p)^{(1-y_i)} \right) \\ &= \sum_{i=1}^n \log \left(p^{y_i} (1-p)^{(1-y_i)} \right) \end{aligned}$$

$$= \sum_{i=1}^n \left(y_i \log p_i + (1 - y_i) \log (1 - p_i) \right)$$

$$\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta} \left(\dots \right)$$

$$= \operatorname{argmin}_{\theta} \left[- \sum_{i=1}^n y_i \log p_i + (1 - y_i) \log (1 - p_i) \right]$$

Cross-Entropy Loss:

$$\frac{1}{n} \sum_{i=1}^n \left(y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i) \right)$$

Loss Function for Logistic Regression

Multi-Class Logistic Regression

Generalization of cross-entropy loss to multi class;

$$-\sum_{i=1}^c y_i \log \hat{y}_i$$

$c = \#$ of classes

For multi-class logistic regression:

$$y = \text{softmax} \left(\frac{\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n}{e^{x_i}} \right)$$

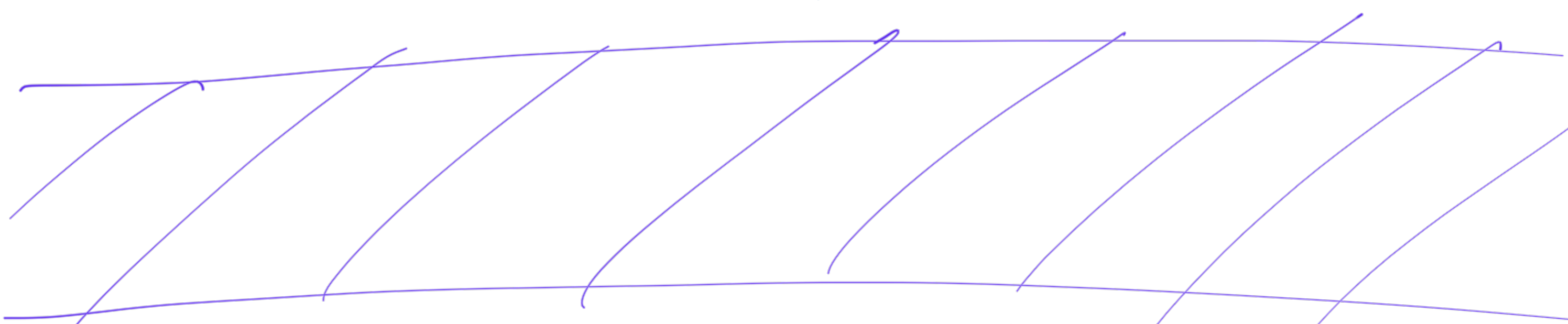
$$\text{Softmax}(X)_i = \frac{e^{X_i}}{\sum_{j=1}^c e^{X_j}}$$

Example:

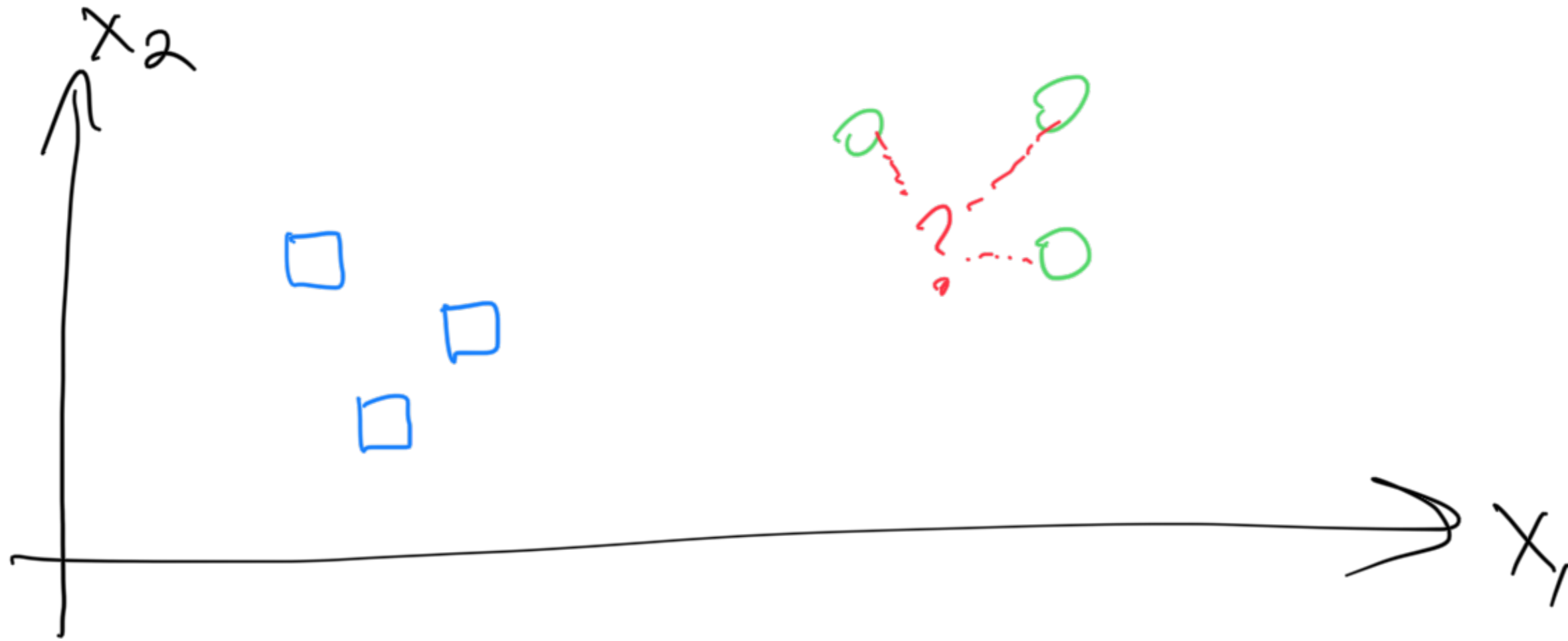
$[-3, 4, 9]$

$$\text{Softmax}(-3) = \frac{e^{-3}}{(e^{-3} + e^4 + e^9)}$$

$$\text{Softmax}(4) = \frac{e^4}{(e^{-3} + e^4 + e^9)}$$

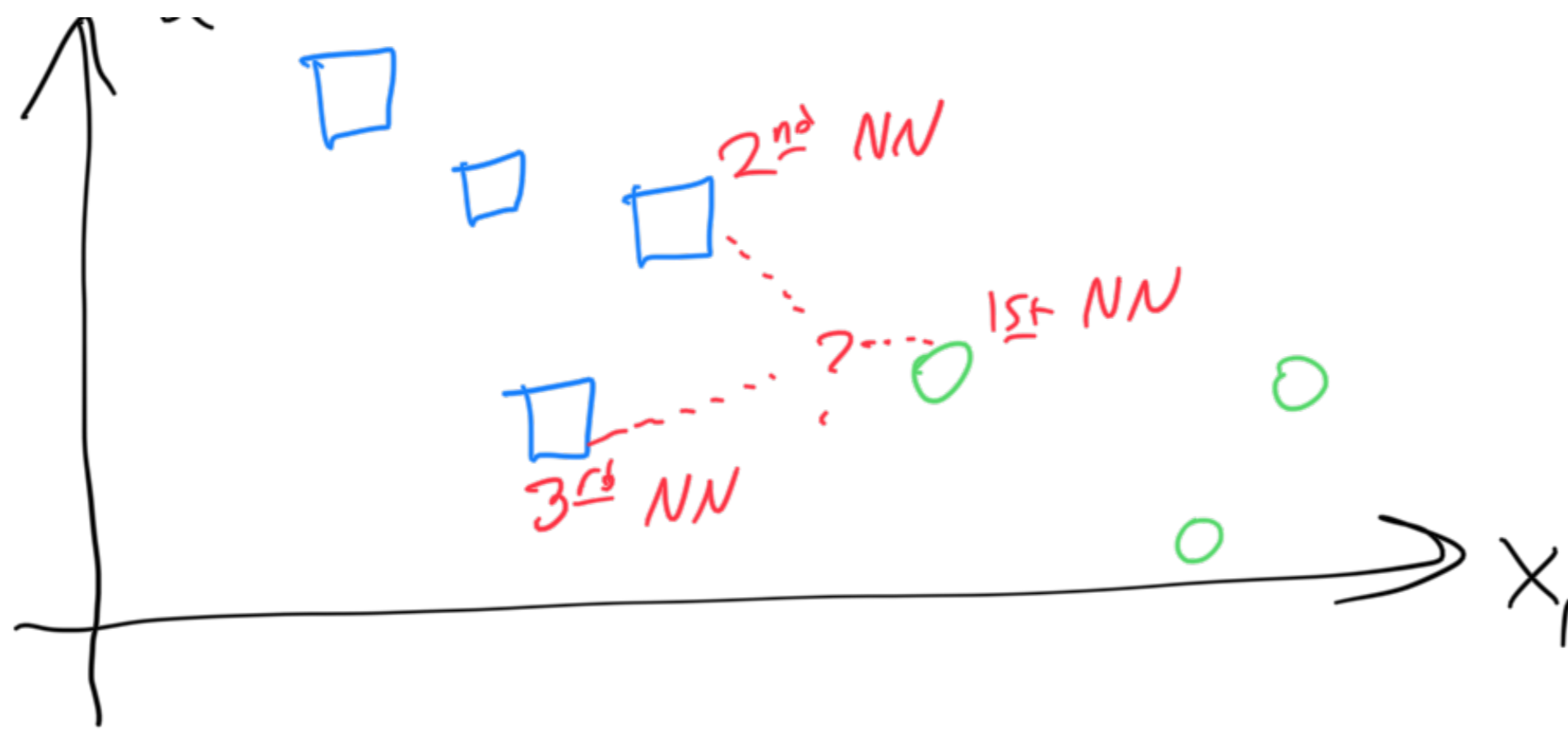


K-Nearest Neighbors (KNN)

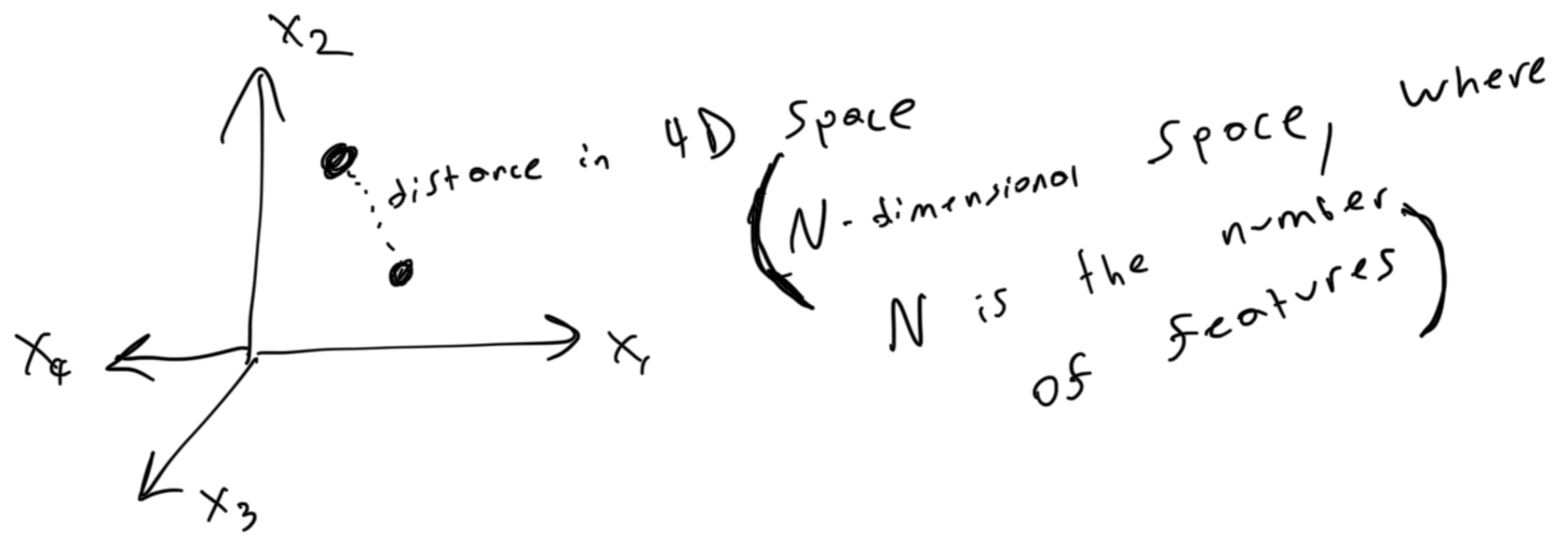


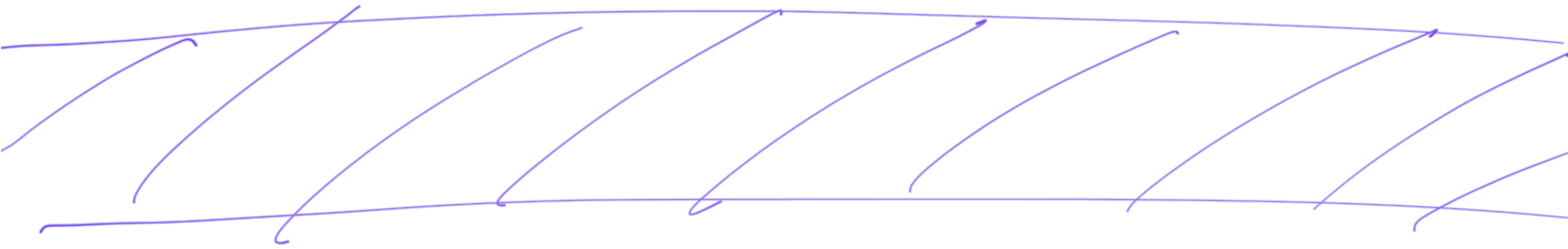
KNN algorithm: Find K nearest neighbors of the point in question and classify as the majority

Choice of K matters: (K is a hyperparameter)



$k = 1$: returns green circle
 $k = 3$: returns blue box





Naive Bayes Classification

Bayes' Rule:

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

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

(Spam Classification)

To Classify:

if $P(y_i = \text{"Spam"} | X_i) > P(y_i = \text{"Not Spam"} | X_i)$
return "Spam"

else:

return "not spam"

To "Train":

$$P(y_i = \text{"Spam"} | X_i) = \frac{P(X_i | y_i = \text{"Spam"}) P(y_i = \text{"Spam"})}{P(X_i)}$$

$$P(y_i = \text{"Spam"}) = \frac{\# \text{ spam messages}}{\# \text{ messages}}$$

$$P(X_i) = \frac{\# \text{ of emails with features } X_i}{\# \text{ messages}}$$

$$P(X_i | y_i = \text{"Spam"}) = \frac{\# \text{ Spam messages with features } X_i}{\# \text{ Spam messages}}$$

The "naïve" in Naïve Bayes:

assume all features X_i are

conditionally independent given label y_i

Example:

$$P(\text{"hello"} = 1, \text{"Vicodin"} = 0, \text{"340"} = 1 | \text{Spam})$$

$$\approx P(\text{"hello"} = 1 | \text{Spam}) \cdot P(\text{"Vicodin"} = 0 | \text{Spam}) \\ \cdot P(\text{"340"} = 1 | \text{Spam})$$

Naive Bayes Example

	x_1	x_2					
$X =$	0	1	,	$g =$			
	1	1					
	0	0					
	1	1					
	1	1					
	0	0					
	1	0					
	1	0					
	1	0					
	1	0					
					0		
					0		
					0		
					0		
					0		
					0		
					0		

x_1 = presence of word x_1
 x_2 = presence of word x_2

Classify $X_? = [1 \ 1]$

$$P(y=0 \mid X = [1 \ 1])$$

$$\propto P(x_1=1 \mid y=0) P(x_2=1 \mid y=0) P(y=0)$$

$$= 1 \cdot 0.25 \cdot 0.4$$

$$= 0.1$$

$$P(y=1 \mid X = [1 \ 1])$$

$$\propto P(x_1=1 \mid y=1) P(x_2=1 \mid y=1) P(y=1)$$

$$= 0.5 \cdot 0.66\bar{6} \cdot 0.6$$

~ ~

$$= 0.2$$

Because $0.2 > 0.1$,
Classify $[1 \ 1]$ as "y=1"