



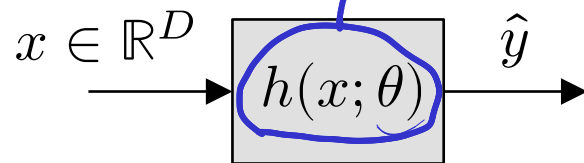
Statistics and Data Science for Engineers

E178 / ME276DS

Classification, Naïve Bayes

Recall

- The prediction problem:



$$\hat{y} = h(x; \theta).$$

- Supervised learning: Given $\mathcal{D} = \{(x_i, y_i)\}_N$

$$\hat{\theta} = \underset{\theta \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^N L(y_i, h(x_i; \theta))$$

- Regression: $y \in \mathbb{R}$

$$L(y, \hat{y}) = L_2(y, \hat{y}) = (y - \hat{y})^2$$

- Linear regression: $h(x, \theta_0, \underline{\theta}_1) = \theta_0 + x^T \underline{\theta}_1 = \theta_0 + \theta_1 x^1 + \dots + \theta_D x^D.$

machine learning.

convex optimization.



explicit solution to linear regression.

$$\hat{\theta} = (X^T X)^{-1} X^T Y \quad \dots \text{explicit} \dots \text{for any } D.$$

Assumption "true system"

$\hat{\theta}$

linear functions of ϵ (noise/uncertainties).

σ^2

→ statistical behavior of $\hat{\theta}$.

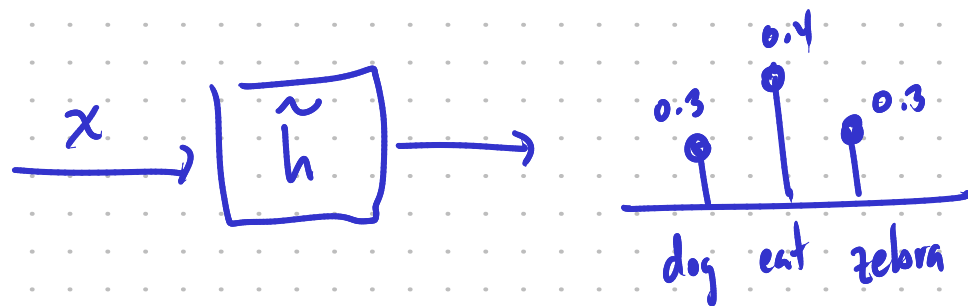
(unbiased and known variance) $\hat{\theta}_0, \hat{\theta}_1, \hat{y}$

statistics.

Var

Classification

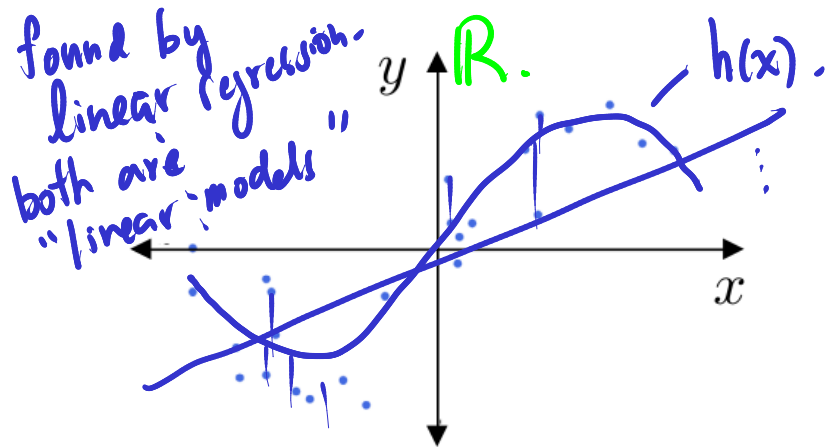
- Sample space: $y \in \{c_1, c_2, \dots, c_K\}$ classes or labels. $\{cat, no\ cat\}.$
 $\{cat, dog, parrot, sign, chair\}.$
- Loss function: ??? $(cat - dog)^2$ $K \dots \# \text{ classes}$ $\{benign, malignant\}$
- Prediction model:
 - Hard classifiers: $\bar{h}(x; \theta)$ returns a *class* $\hat{y} \in \{c_1, c_2, \dots, c_K\}$
 - Soft classifiers: $\tilde{h}(x; \theta)$ returns a *distribution* $P(\hat{y} \mid X = x)$ over classes.



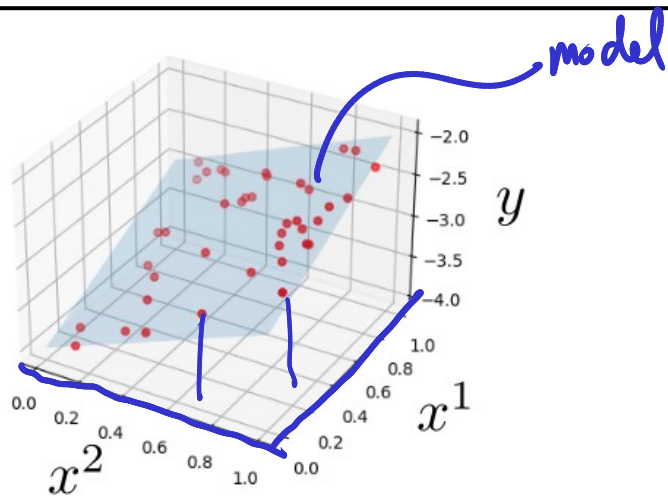
$$\bar{h}(x) = \underset{c}{\operatorname{argmax}} \tilde{h}(x).$$

Regression

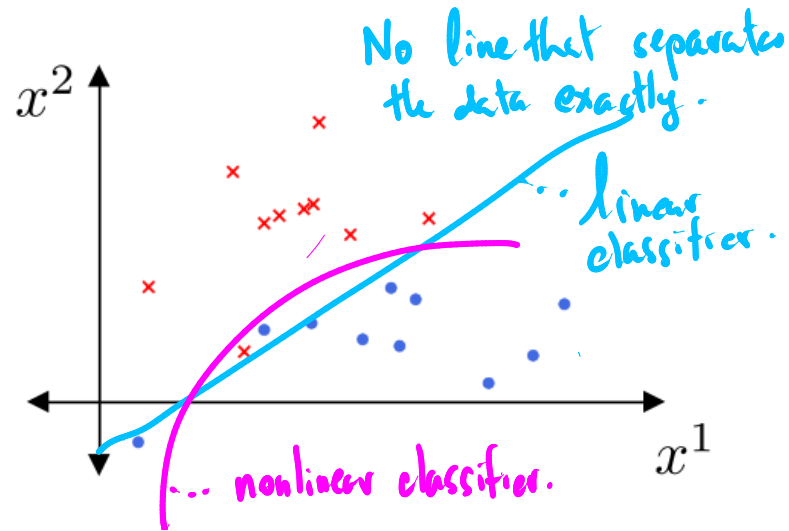
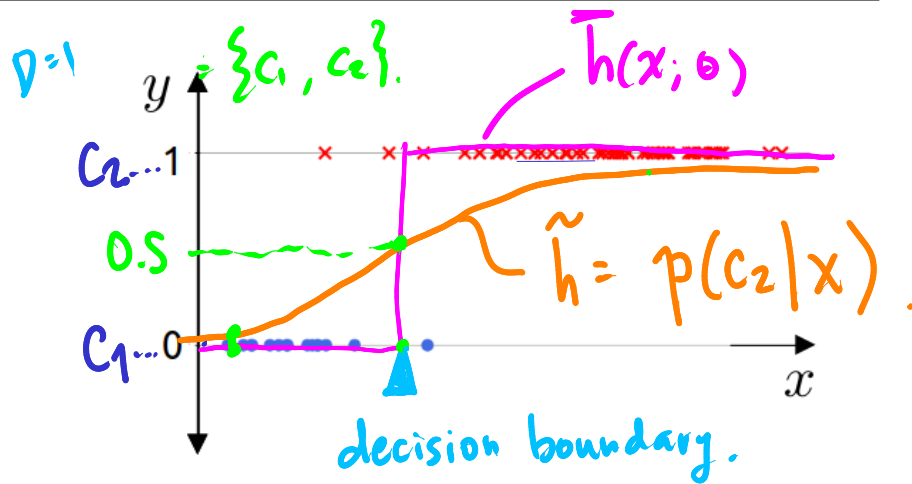
$D=1$

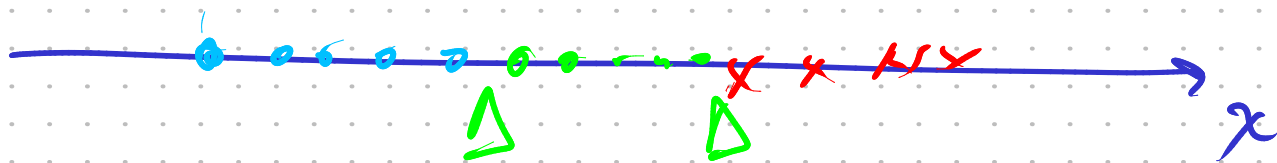


$D=2$

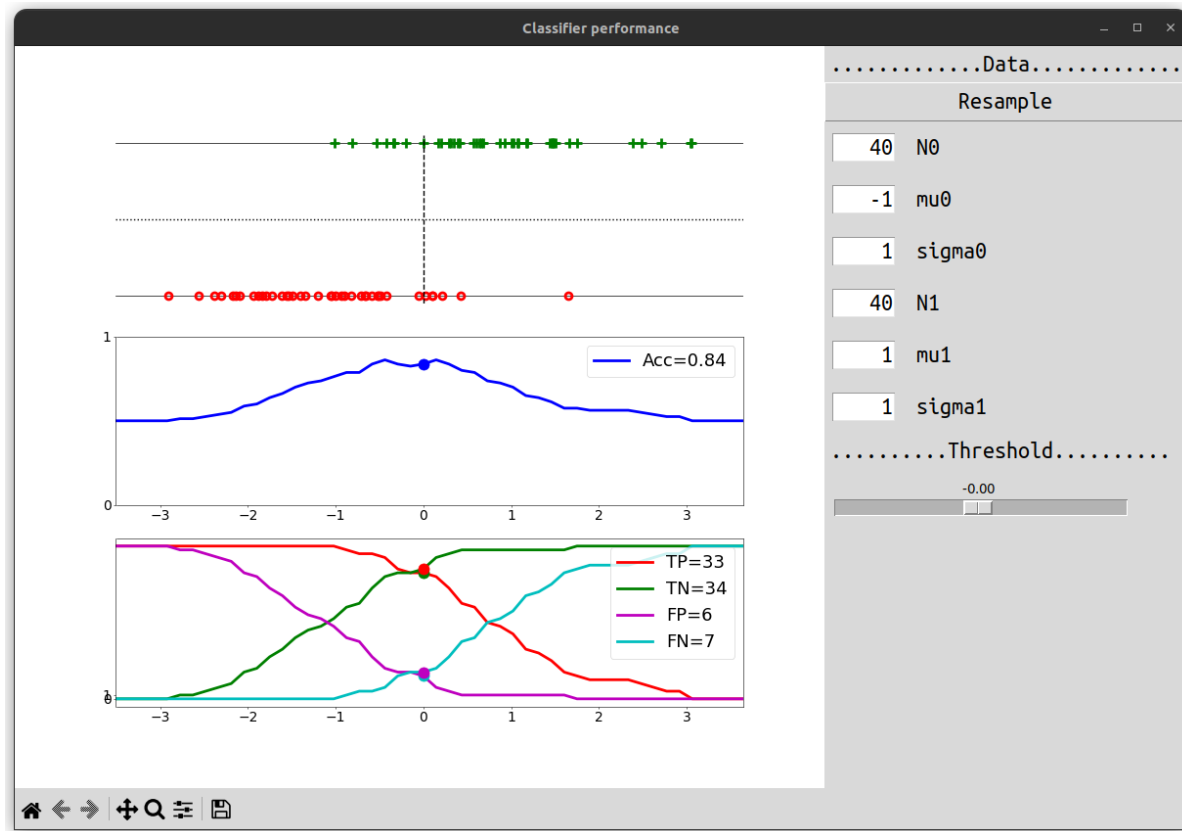


$K=2$ Two-class classification





Demo: 1D, 2-class classification



positives

negative.

False
negatives

True positives
TP

True negatives.

False positives.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} = \frac{TP + TN}{N}$$

$h(x) = \text{negative}$

$h(x) = \text{positive}$

Accuracy is good when the problem is "symmetric"

- I have equal preference for FP and FN.

- The number of data points in each category is approx the same.
(balanced dataset).

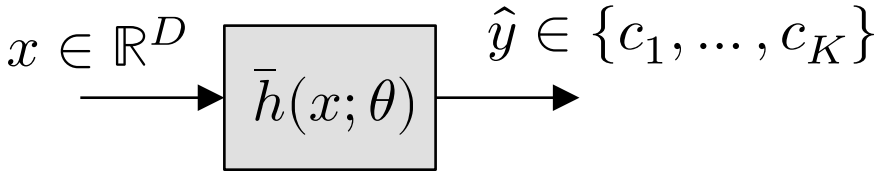
If "non-symmetric" then:

→ precision, recall, balanced accuracy, TPR, TNR

→ ROC curve

→ AUC.

Accuracy as a loss function

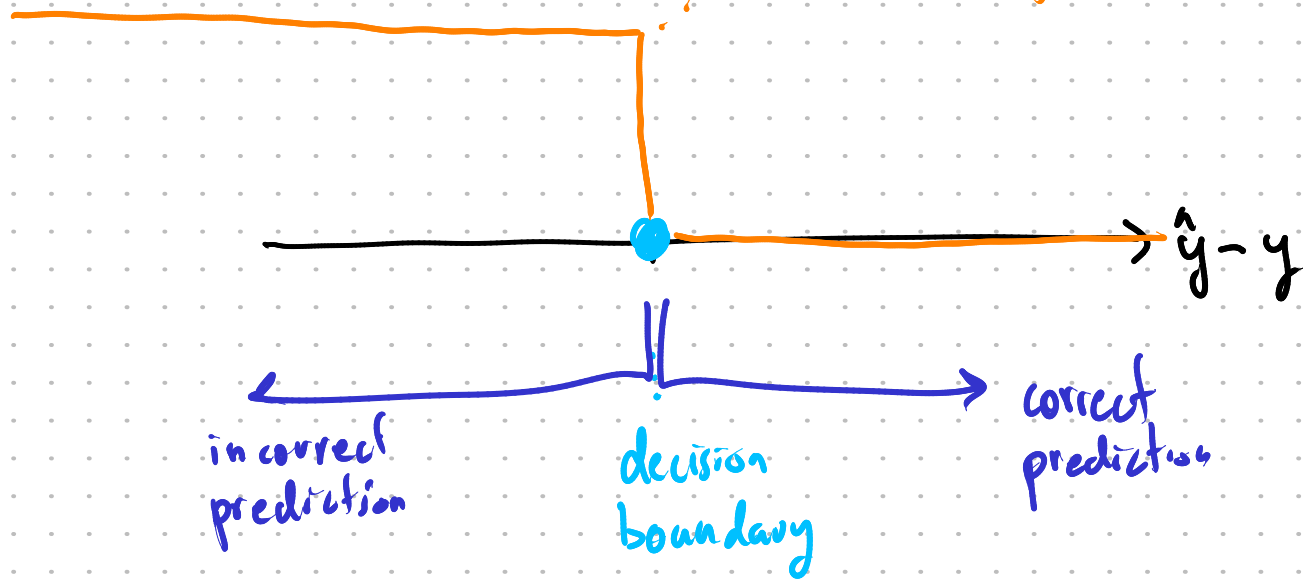
$$L_{01}(y, \hat{y}) = \begin{cases} 1 & y \neq \hat{y} \\ 0 & y = \hat{y} \end{cases}$$


The diagram shows an input $x \in \mathbb{R}^D$ entering a box labeled $\bar{h}(x; \theta)$. An arrow points from the box to the output $\hat{y} \in \{c_1, \dots, c_K\}$.

Accuracy is maximized by solving.

$$\text{minimize } \sum_{i=1}^n L_{01}(y_i, h(x_i; \theta)).$$

L01 loss function. \Rightarrow no gradient descent.



- Cannot solve optimization problem directly with L_1

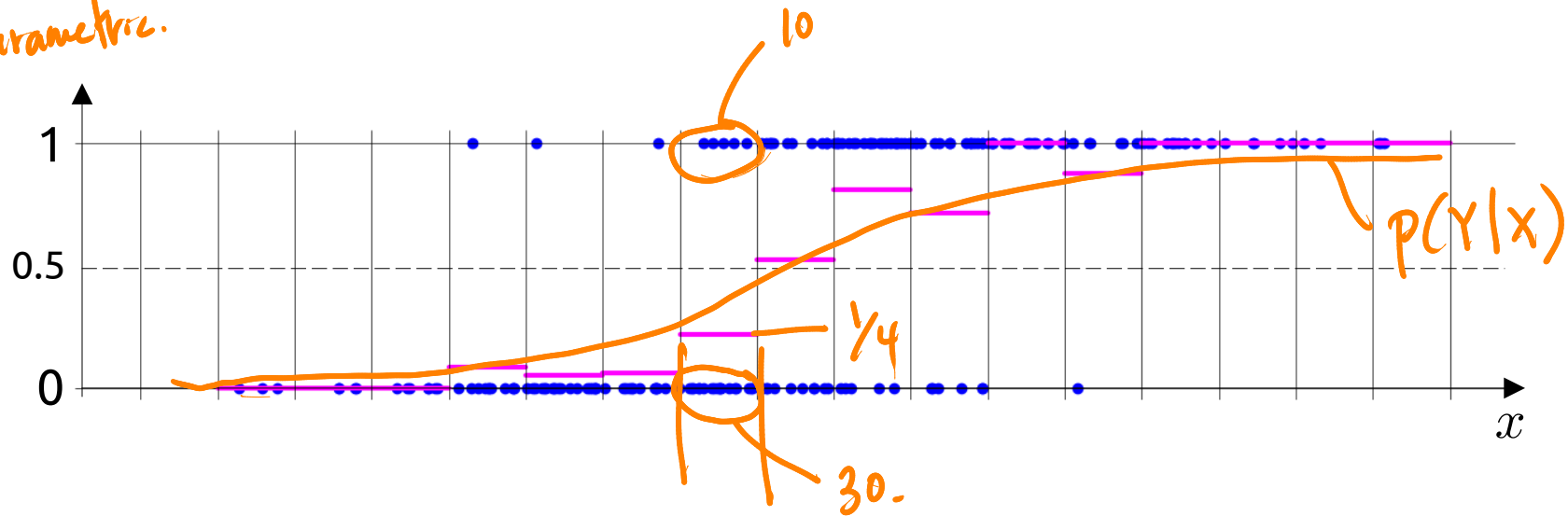
- Assume that we know $p(Y|X)$.



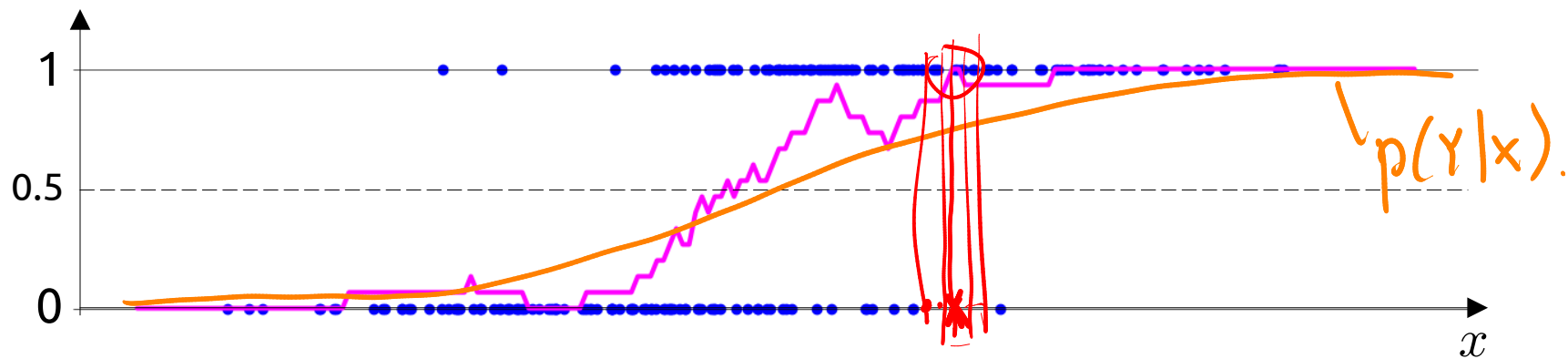
$$\hat{y} = \underset{c}{\operatorname{argmax}} p(Y|X=x) \Rightarrow \text{minimize } E[L_1] \\ \Rightarrow \text{maximize accuracy.}$$

- Approximate $p(Y|X=x)$.

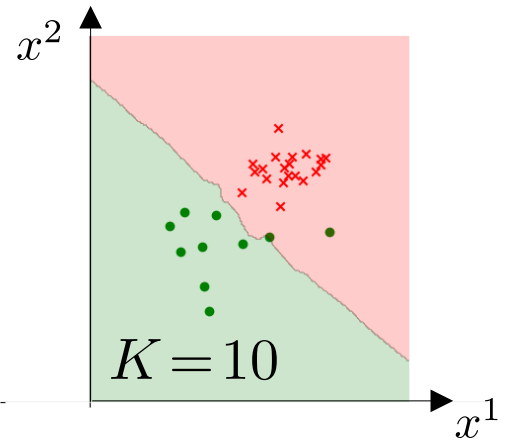
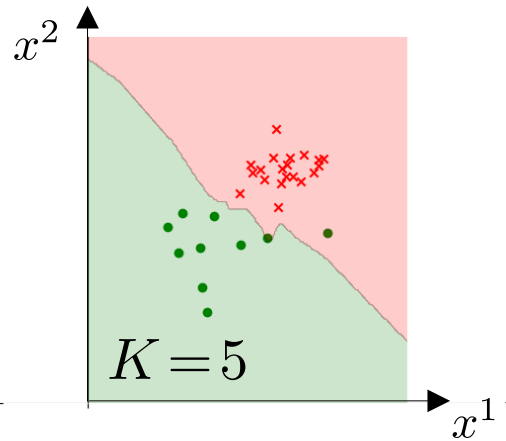
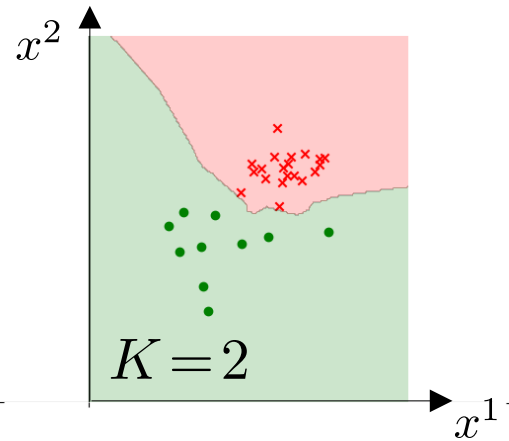
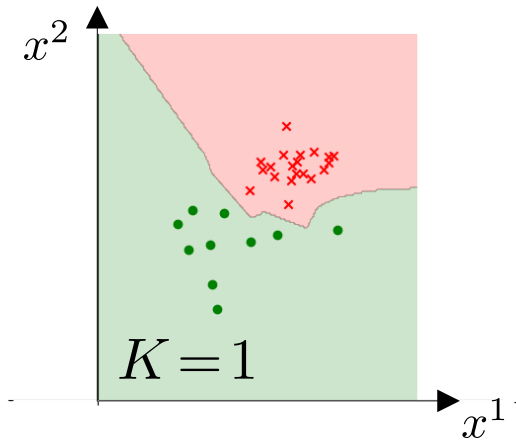
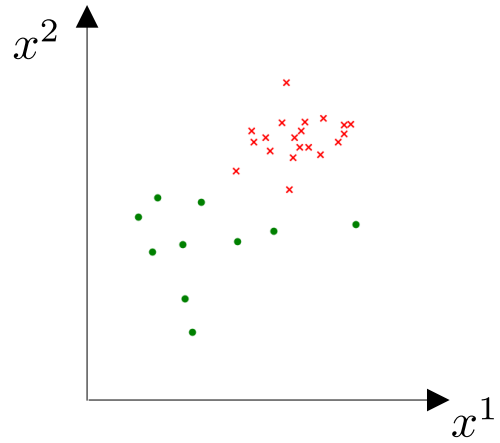
parametric.
K-bins



KNN



Example: Classifying two-class / 2D data with KNN



Applying Bayes' rule

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{c \in \{c_1, \dots, c_K\}} P(Y = c \mid X = x) && \textcircled{\text{I}} \\ &= \operatorname{argmax}_{c \in \{c_1, \dots, c_K\}} \frac{p(X = x \mid Y = c) P(Y = c)}{p(X = x)} && \text{Bayes' rule.} \\ &= \operatorname{argmax}_{c \in \{c_1, \dots, c_K\}} p(X = x \mid Y = c) P(Y = c) && \textcircled{\text{II}}\end{aligned}$$

posterior belief.

likelihood of the input x given the class.

prior belief.

① Estimate a discrete distribution for every input $x \in \mathbb{R}^D$.

② Estimate a continuous distribution for every class + 1 discrete distr.

$$\hat{y} = \underset{c \in \{c_1, \dots, c_K\}}{\operatorname{argmax}} \quad \underbrace{P(Y=c)}_{\text{prior}} \underbrace{p(X^1=x^1, X^2=x^2 | Y=c)}_{\text{likelihood}}$$

$$P(Y=X) = 20/30$$

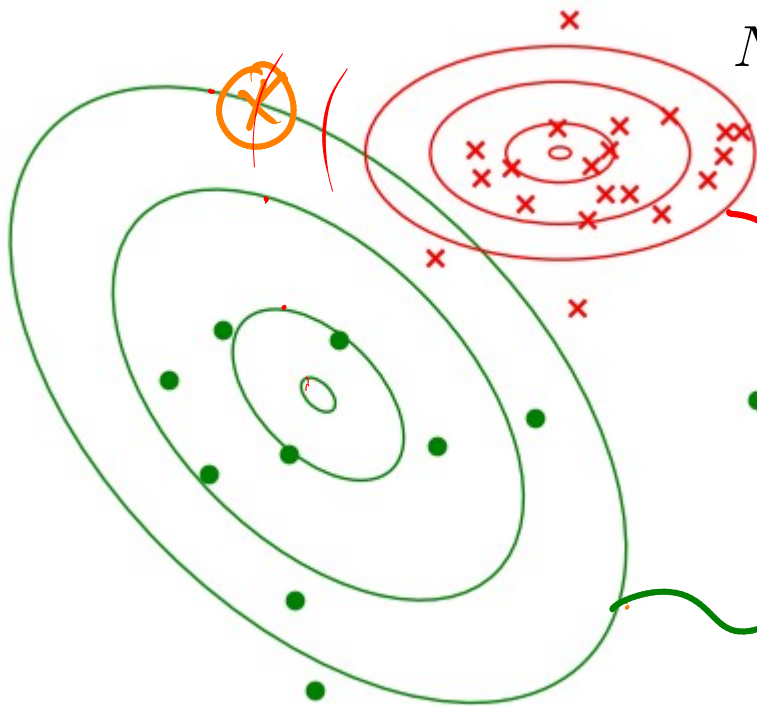
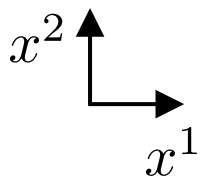
$$P(Y=0) = 10/30$$

$$N^0 = 10$$

$$N^X = 20$$

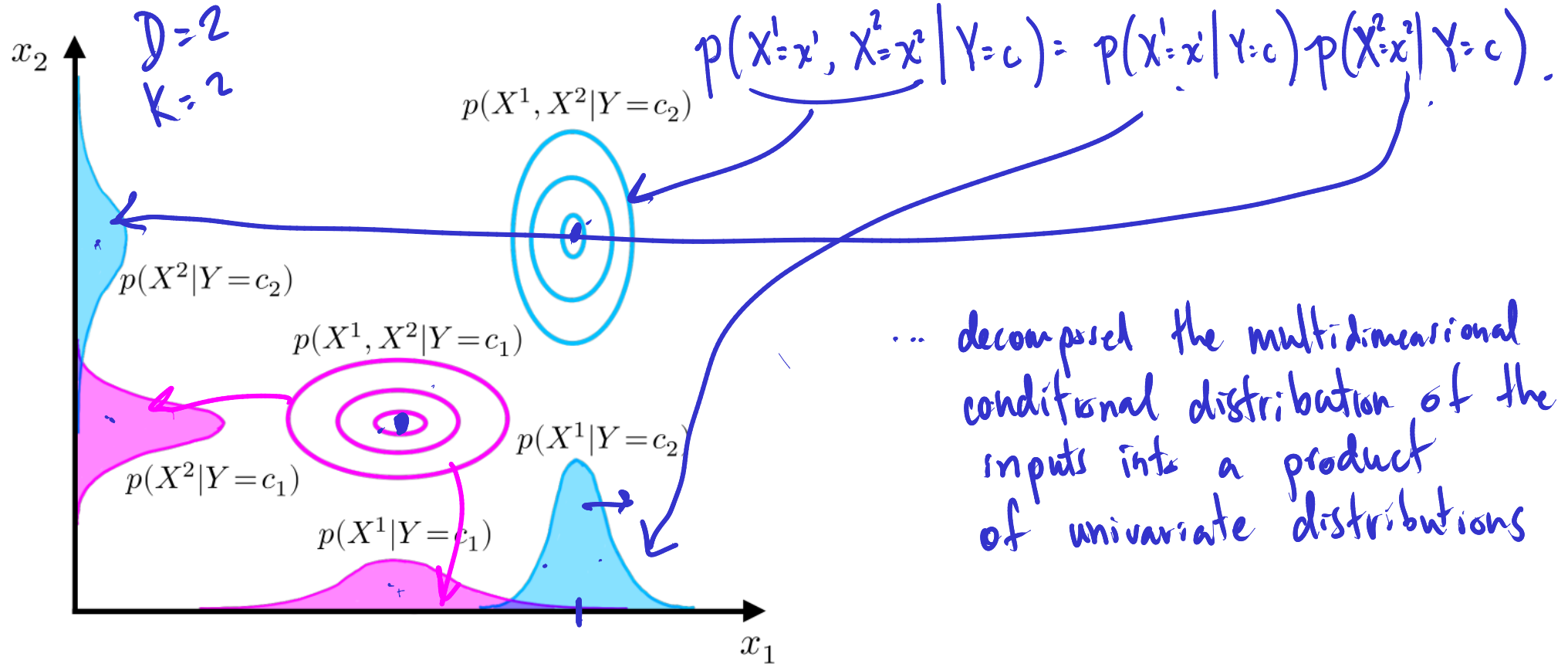
$$p(X=x | Y=x)$$

$$p(X=x | Y=0)$$



Naïve Bayes

Assumption: The components of the input are independent given the class.



Original number of parameters.

K 2D Gaussians $\begin{cases} D \text{ parameters for the mean} \\ O(D^2) \text{ parameter for the covariance.} \end{cases}$
 $\rightarrow K \cdot D^2$

Number of parameters after assumption.

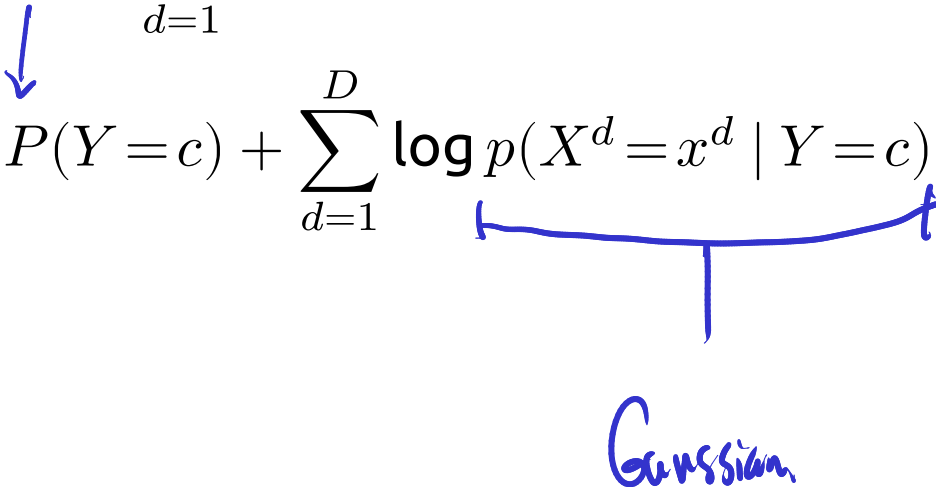
$K \cdot D \cdot 2.$

$$O(KD^2) \longrightarrow O(KD).$$

Naïve Bayes

Assumption: The components of the input are independent given the class.

$$\begin{aligned}\hat{y} &= \operatorname{argmax}_{c \in \{c_1, \dots, c_K\}} P(Y=c) \prod_{d=1}^D p(X^d = x^d \mid Y=c) \\ &= \operatorname{argmax}_{c \in \{c_1, \dots, c_K\}} \log P(Y=c) + \sum_{d=1}^D \log p(X^d = x^d \mid Y=c)\end{aligned}$$



Gaussian

Gaussian Naïve Bayes

Assumption: The individual class-conditioned inputs are Gaussian

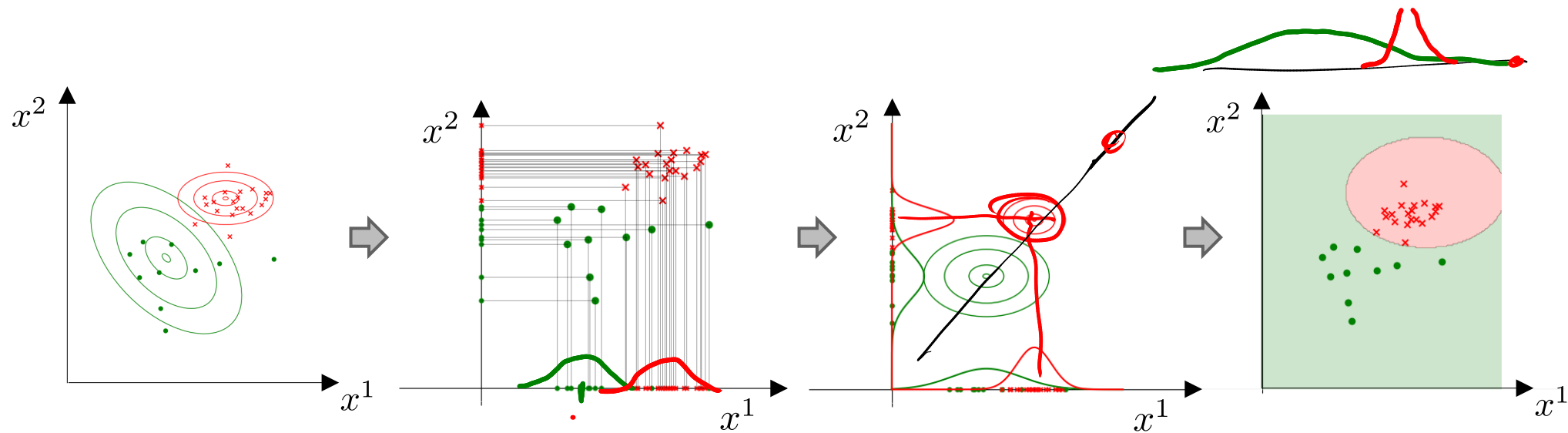
$$X^d = x^d | Y = c \sim \mathcal{N}(\mu_{d,c}, \sigma_{d,c}^2) \quad \forall d, c$$

Training: Compute point estimates of the $\mu_{d,c}$'s and $\sigma_{d,c}^2$'s (and of the $P(Y=c)$'s)

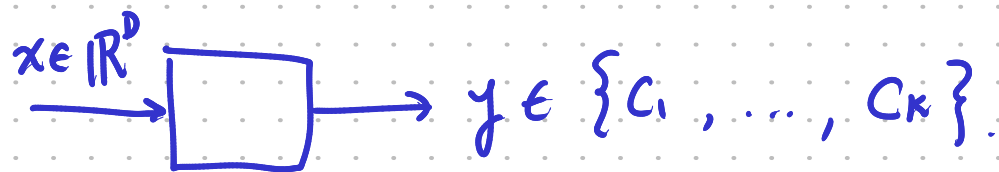
Prediction: Choose the class that maximizes the posterior probability:

$$\hat{y} = \underset{c \in \{c_1, \dots, c_K\}}{\operatorname{argmax}} \quad \log P(Y=c) + \sum_{d=1}^D \log p(X^d = x^d | Y=c)$$

Example: Classifying two-class / 2D data with Gaussian NB



Recap:



- Maximize accuracy.



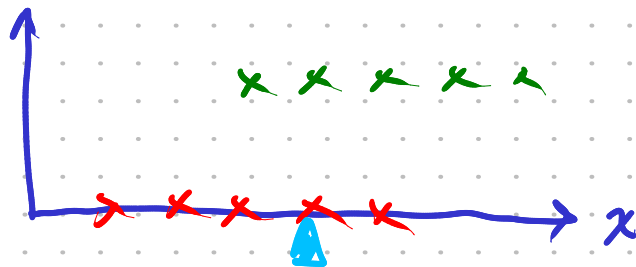
- Minimize L_{01} ... problem: L_{01} cannot be minimized easily.



- Maximize $P(Y=c|X=x)$ Don't know $p(Y|x)$.

\Downarrow Bayes' rule.

Maximize $P(X=x|Y=c) P(Y=c)$ approximate these probabilities



↓ Naïve Bayes assumption.

$$p(X=x|Y=c) = \prod_{d=1}^D p(x^d=x^d|Y=c).$$

Naïve Bayes

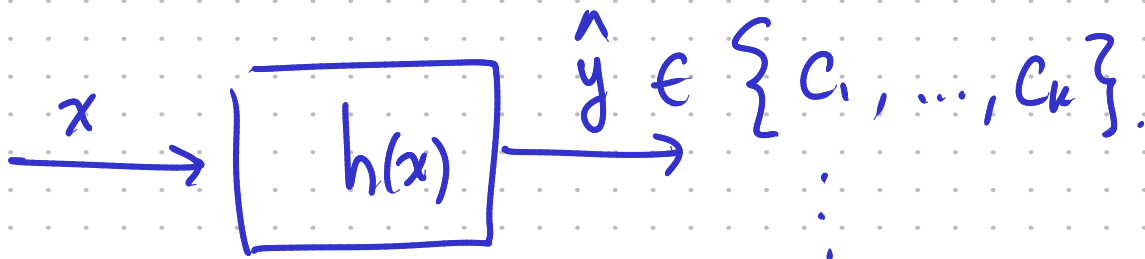
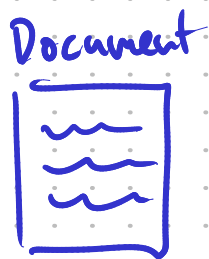
Inputs are label based.

⋮
text classification.

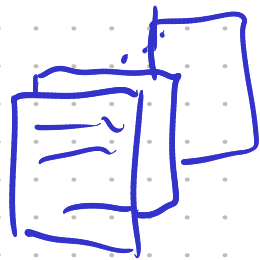
Assume
 $p(x^d=x^d|Y=c)$
are Gaussian

Gaussian Naïve Bayes.

Text classification



Training data: "corpus"



$\{ \text{children's book, romance novel, sci-fi story} \dots \}$

$\{ \text{positive, negative} \}$

How do we encode the document into D inputs. : Bag-of-words.

Naïve Bayes with labeled inputs

$$x^1 \in \{a, b, c\}$$

$$x^2 \in \{\alpha, \beta, \gamma\}$$

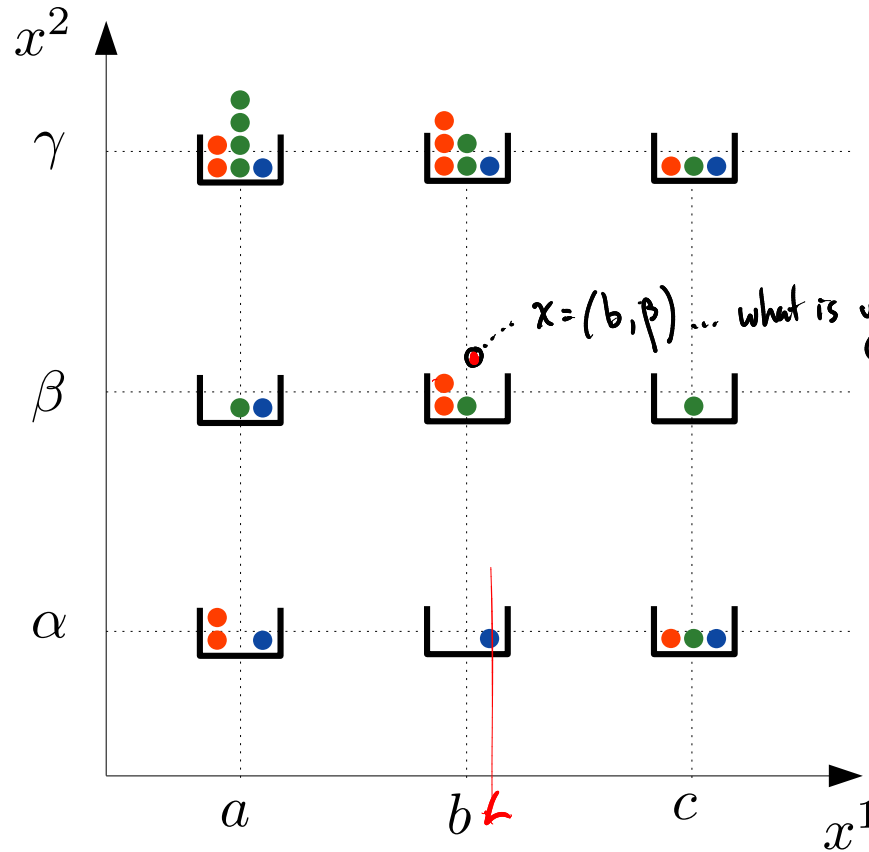
$$y \in \{\text{red dot}, \text{green dot}, \text{blue dot}\}$$

$$P(Y=\text{red dot}) \approx \frac{\# \text{ red dots}}{N} = 11/29 = \hat{P}_R$$

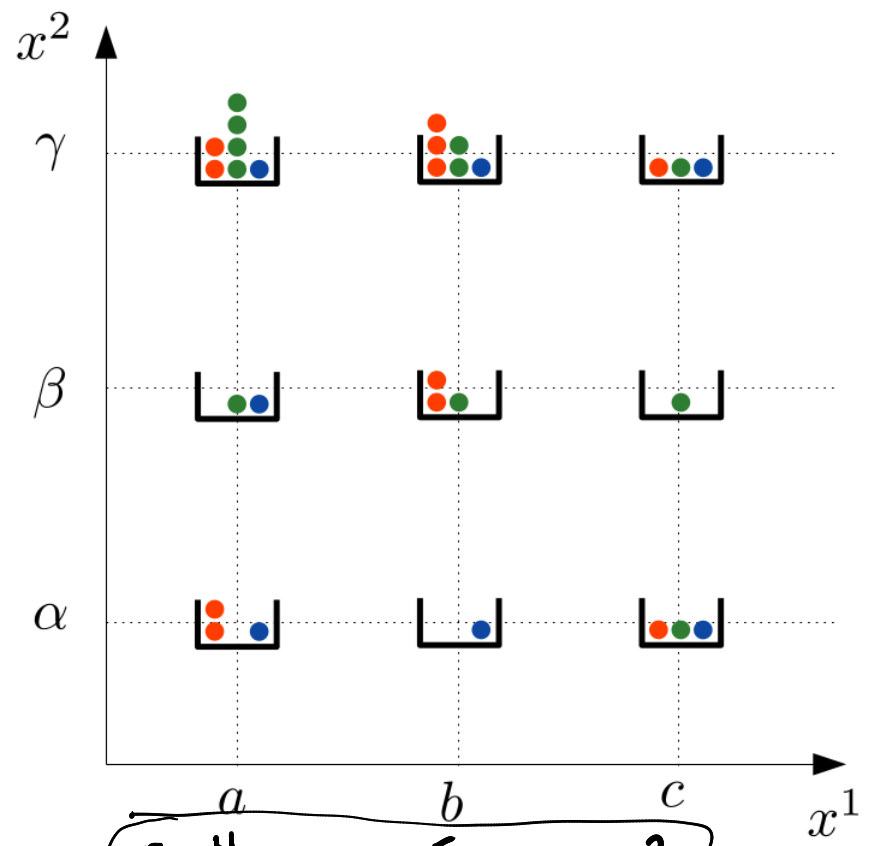
$$P(Y=\text{green dot}) = 11/29 = \hat{P}_G$$

$$P(Y=\text{blue dot}) = 7/29 = \hat{P}_B$$

P



r	g	b
6	7	3
2	3	1
3	1	3



r	4	5	2
g	5	3	3
b	3	2	2
	12	10	7

1

$$\hat{y} = \underset{c}{\operatorname{argmax}} \underbrace{\log P(Y=c)} + \sum_{d=1}^D \log \underbrace{p(X^d = x^d | Y=c)} \quad \star$$

$$p(X^1 = b | Y = \bullet) = 5/11$$

$$p(X^2 = \beta | Y = \bullet) = 2/11$$

$$\hat{y} = \underset{c}{\operatorname{argmax}} \left(\underbrace{\log \frac{11}{29} + \log \frac{5}{11} + \log \frac{2}{11}}_{\text{red}}, \underbrace{\hspace{10em}}_{\text{green}}, \underbrace{\hspace{10em}}_{\text{blue}} \right)$$

Train the model: Use $\mathcal{D}_{\text{train}}$ to compute

$$\hat{p}(Y=c)$$

$$\hat{p}(X^d = x^d | Y=c)$$

Evaluate/Prediction: evaluating $\text{argmax}(\log \dots)$ ★

Naïve Bayes with labeled inputs

$$\hat{y} = \operatorname{argmax}_{c \in \{\text{red}, \text{green}, \text{blue}\}} \log P(Y=c) + \sum_{d=1}^D \log p(X^d = x^d \mid Y=c)$$

$$= \operatorname{argmax}_{c \in \{\text{red}, \text{green}, \text{blue}\}} \left(\begin{array}{l} \log \frac{11}{29} + \log \frac{5}{11} + \log \frac{2}{11}, \quad \log \frac{11}{29} + \log \frac{3}{11} + \log \frac{3}{11}, \quad \log \frac{7}{29} + \log \frac{2}{7} + \log \frac{1}{7} \end{array} \right)$$

red green blue.

$$= \operatorname{argmax}_c \left(\cancel{\frac{11}{29}} \frac{5}{11} \frac{2}{11}, \quad \cancel{\frac{11}{29}} \frac{3}{11} \frac{3}{11}, \quad \cancel{\frac{7}{29}} \frac{2}{7} \frac{1}{7} \right)$$

$$= \operatorname{argmax}_c \left(\frac{10}{11}, \frac{9}{11}, \frac{2}{7} \right) = \text{red}$$

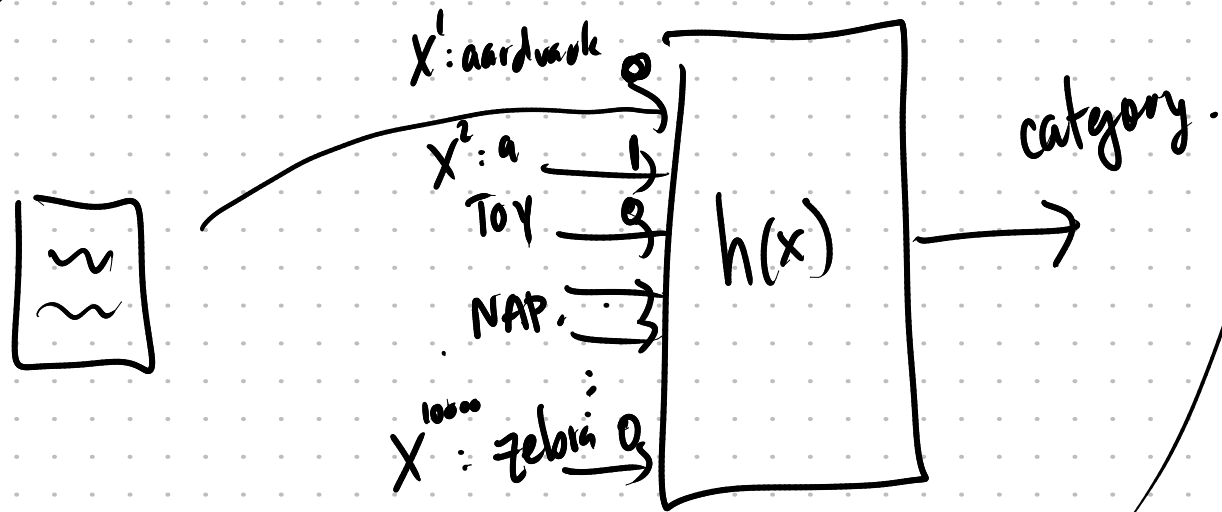
Example: Text classification

Problem

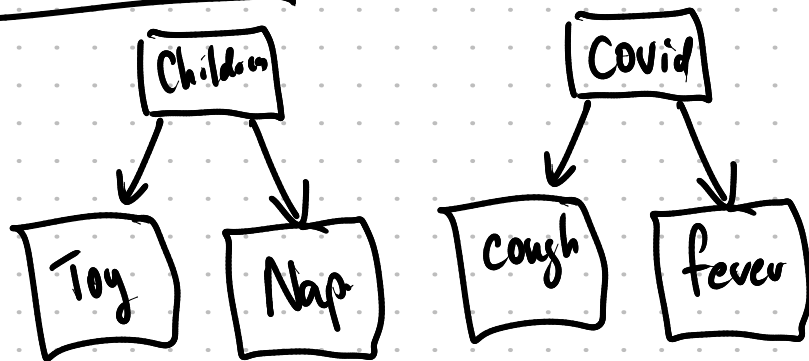
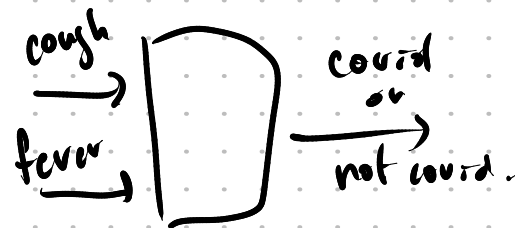
- A sample x is a *document* (article, blog entry, tweet, review, etc).
- Categories $\{c_1, c_2, \dots, c_K\}$ (e.g. positive/negative, genre, etc).
- Training data (a *corpus*): $\mathcal{D} = \{(x_i, y_i)\}_N$

Input encoding: Vocabulary of D words.

X^1	= a is present	
X^2	= aardvark is present	(all binary)
\vdots		
$X^{10,000}$	= zebra is present	



Aside on the Naive Bayes assumption. that the inputs are independent given the output:



Bag of words:

Handwritten annotations:

- A bracket above the word columns points to the label "y: target".
- A handwritten "1" is circled in the first row, first column.
- A handwritten "1" is circled in the eighth row, second column.

	a	abacus	abandon	amazing	...	free	...	class
1	1	0	0	0	...	0	...	c1
1	1	0	0	0	...	0	...	c1
1	1	0	1	0	...	0	...	c3
1	1	0	0	0	...	0	...	c1
1	1	0	0	0	...	0	...	c2
1	1	0	0	0	...	0	...	c3
1	1	1	0	0	...	0	...	c1
1	1	0	0	0	...	0	...	c2
1	1	0	0	0	...	0	...	c1
1	1	0	0	1	...	0	...	c2

$$\hat{y} = \operatorname{argmax}_{c \in \{c_1, c_2, c_3\}} \underbrace{\log P(Y=c)}_{\text{... vocabulary.}} + \sum_{d=1}^D \underbrace{\log p(X^d = x^d | Y=c)}_{\text{...}}$$

$$P(Y=c) \approx \frac{\text{\# docs of class } c}{\text{\# docs}} = \frac{N_c}{N} = \hat{p}_c$$

... word d is present.

$$\rightarrow P(\underbrace{X^d=1}_{\text{word } d \text{ is present}} | Y=c) \approx \frac{\text{\# docs of class } c \text{ with word } d}{\text{\# docs of class } c} = \boxed{\frac{N_{d,c}}{N_c} = \hat{p}_{d,c}}$$

word d is not present

$$P(\underbrace{X^d=0}_{\text{word } d \text{ is not present}} | Y=c) \approx 1 - \hat{p}_{d,c}$$

Problems:

- 1) Empty class: $N_c = 0 \Rightarrow N_{d,c} = 0 \quad \forall d \Rightarrow \hat{p}_{d,c}$ is undefined $\forall d$
- 2) New word: $N_{d,c}$ is undefined

Solution: Laplace smoothing

$$\hat{p}_{d,c} = \frac{N_{d,c} + \alpha}{N_c + \alpha K}$$

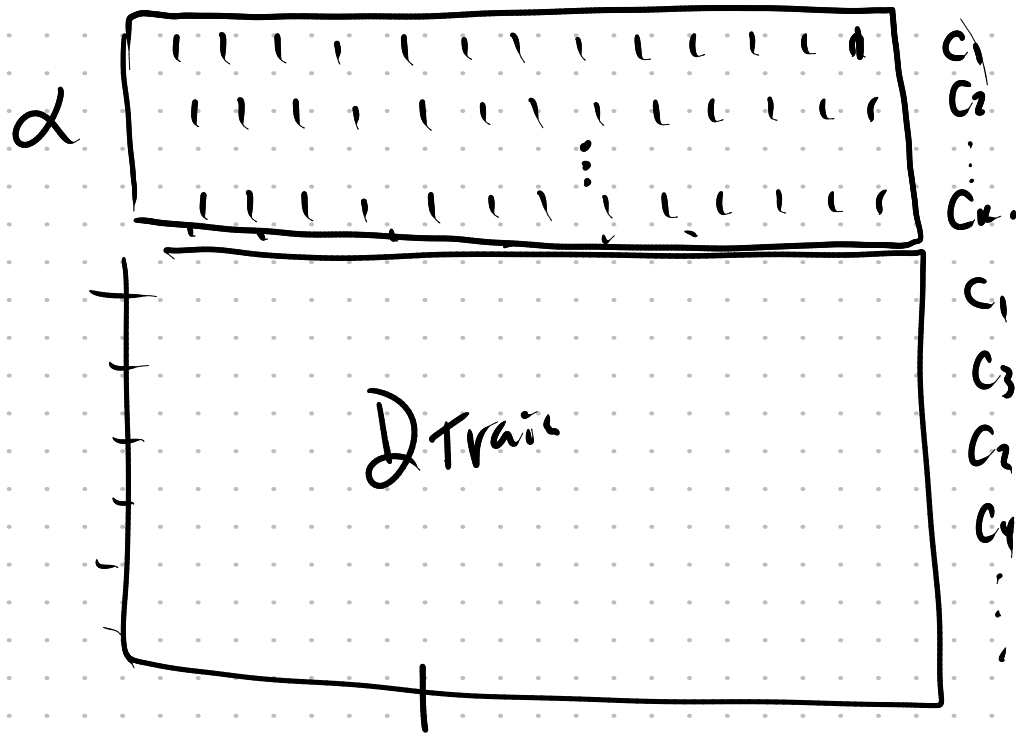
hyperparameter.

\vdots

$$\alpha > 0$$

$$N_{d,c} \rightarrow N_{d,c} + 1 \cdot \alpha$$

$$N_c \rightarrow N_c + \underbrace{K}_K \cdot \alpha$$



$N_{d,c}$ are increased
by 1.

All N_c 's are increased by K

Predict the class of this new document:

a	abacus	abandon	amazing	...	free
1	0	1	0	...	0	...

$$\hat{y} = \operatorname{argmax}_c \log P(Y=c) + \sum_{d=1}^D \log P(X^d = x^d \mid Y=c)$$

$$\approx \operatorname{argmax}_c \log \hat{p}_c + \sum_{d: x^d=1} \log \hat{p}_{d,c} + \sum_{d: x^d=0} \log(1 - \hat{p}_{d,c})$$

\vdots
word in the document

\vdots
words not in the document.

Sorted table:

a	abacus	abandon	amazing	...	free	class
1	0	0	0	...	0	...	c1
1	0	0	0	...	0	...	c1
1	0	0	0	...	0	...	c1
1	1	0	0	...	0	...	c1
1	0	0	0	...	0	...	c1
1	0	0	0	...	0	...	c2
1	0	0	0	...	0	...	c2
1	0	0	1	...	0	...	c2
1	0	1	0	...	0	...	c3
1	0	0	0	...	0	...	c3