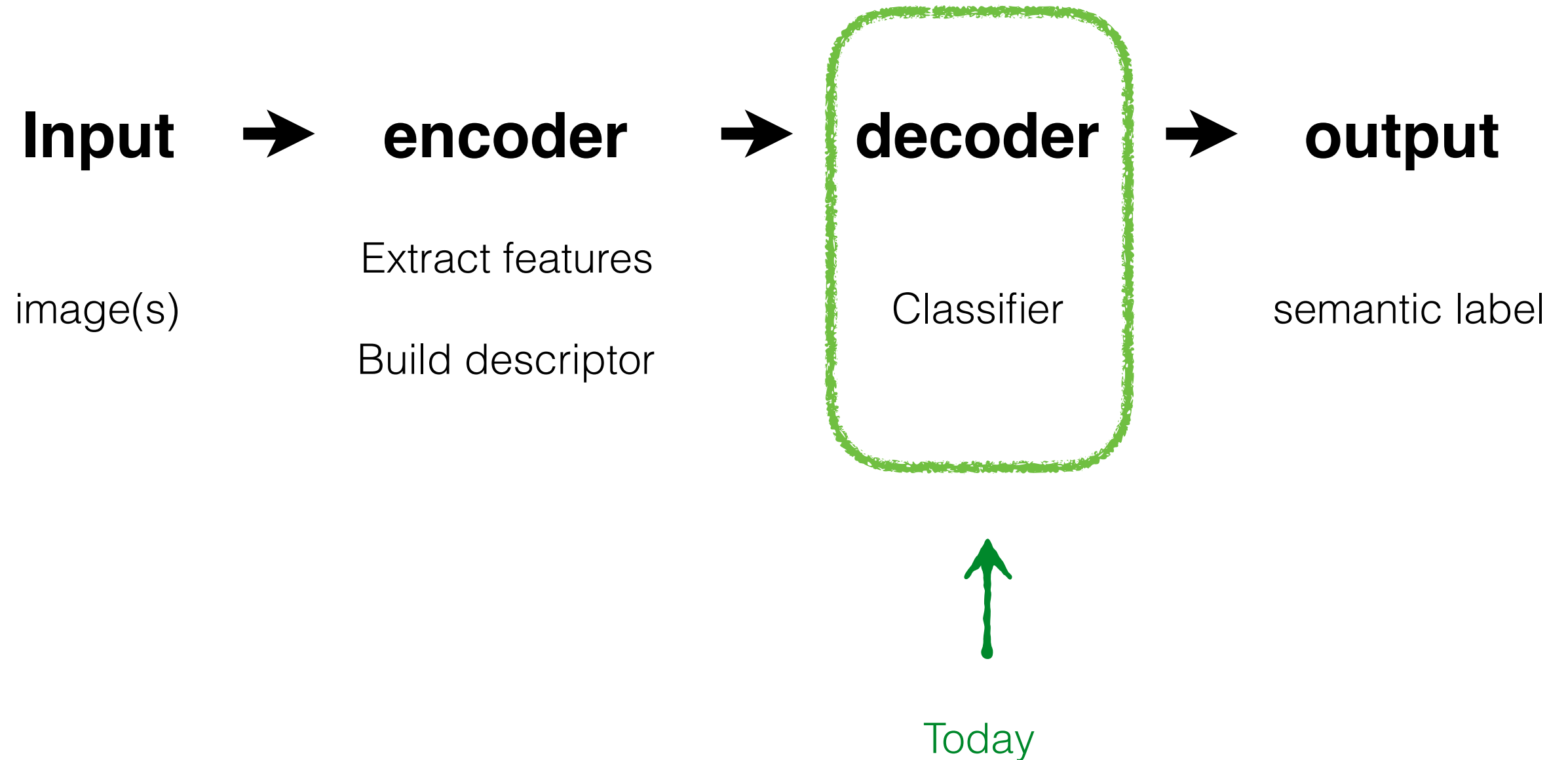
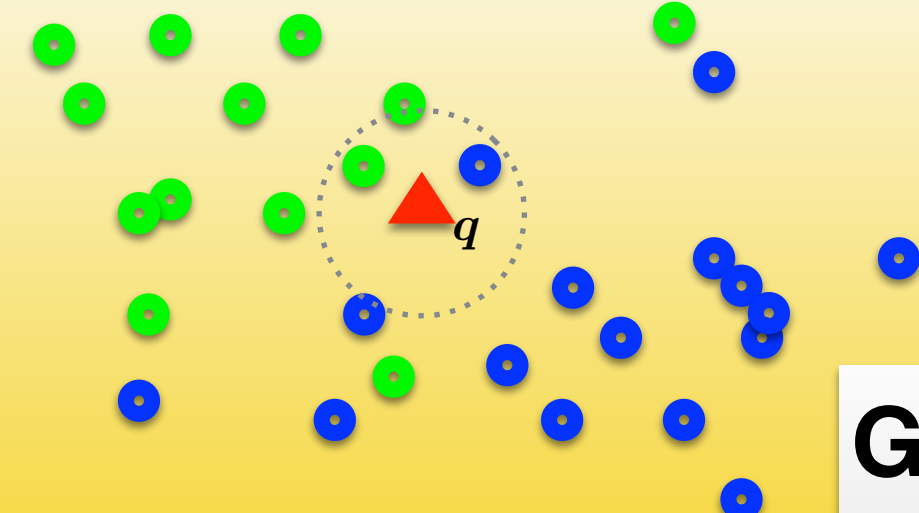


‘Classical’

Image Classification Pipeline



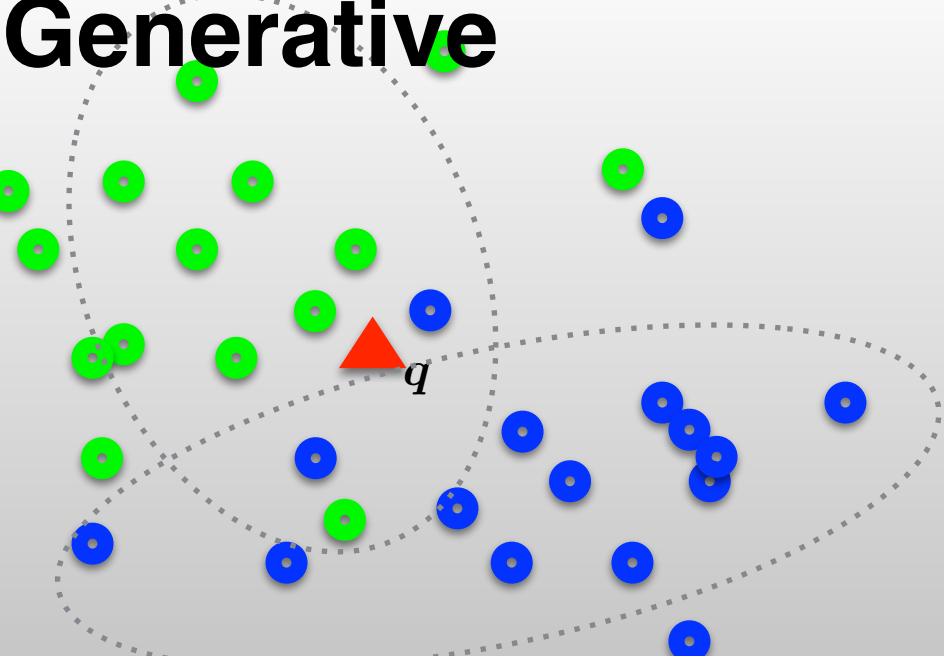
Non-parametric



Nearest Neighbor

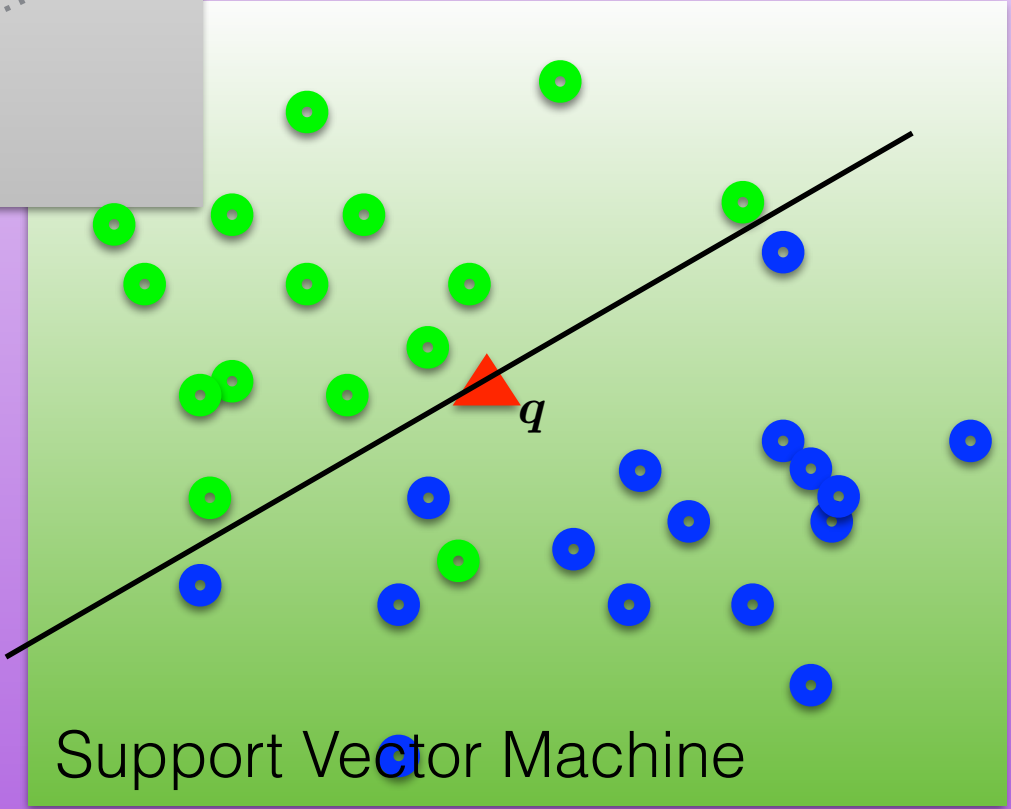
Parametric

Generative



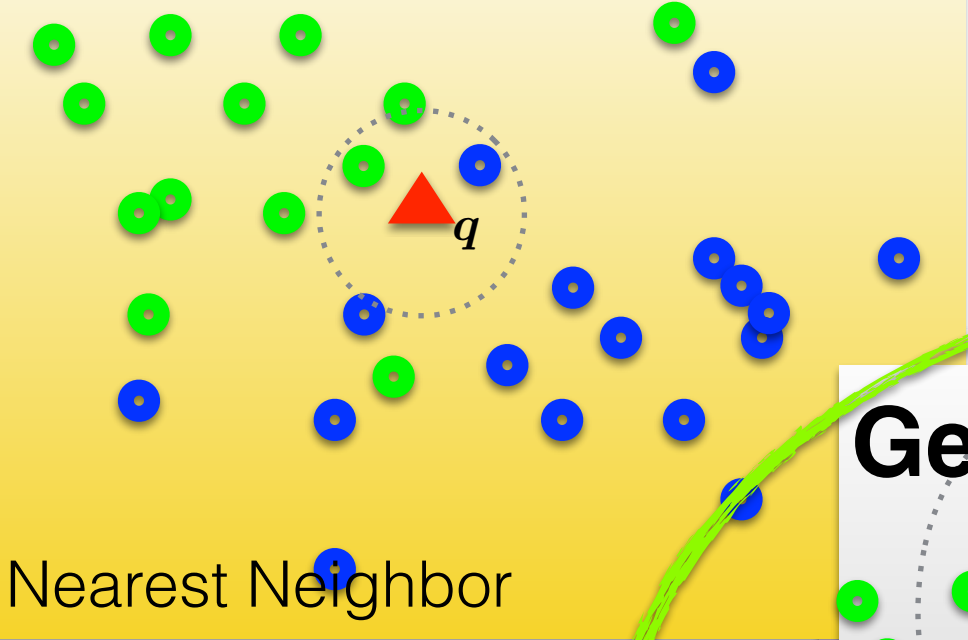
Naive Bayes

Discriminative

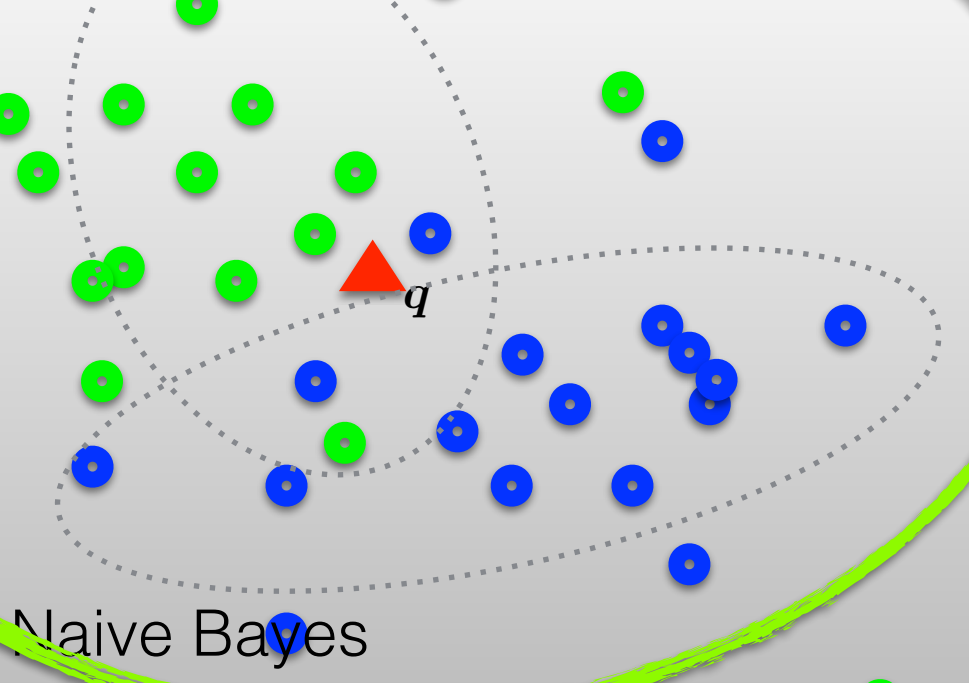


Support Vector Machine

Non-parametric

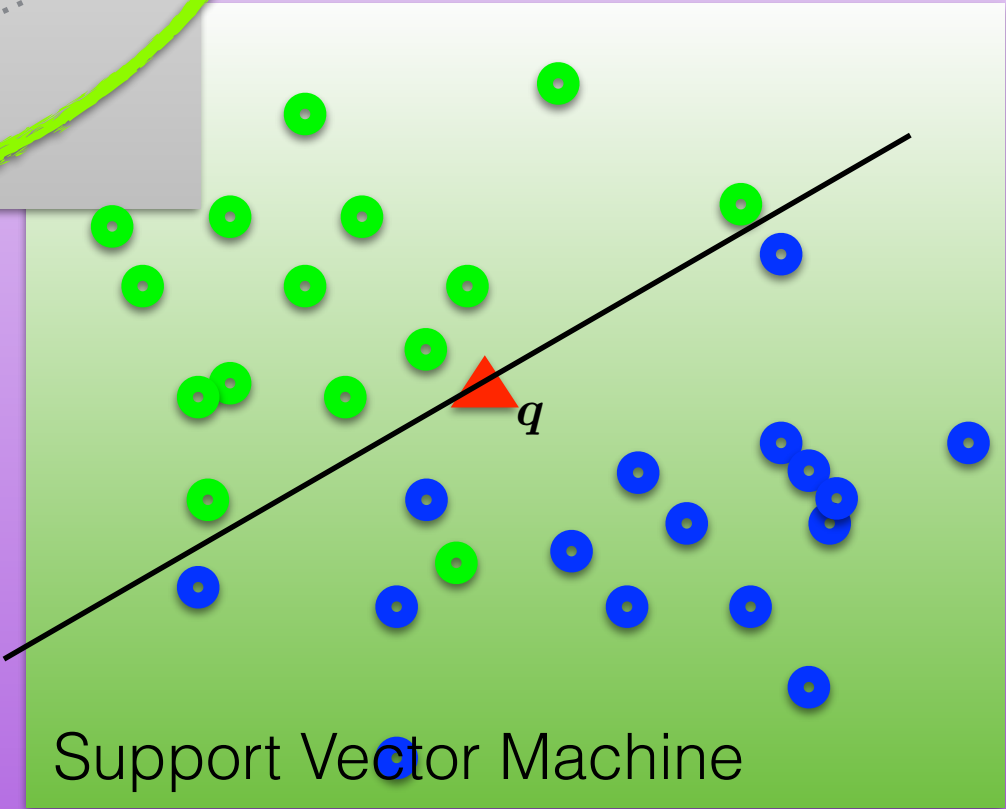


Generative

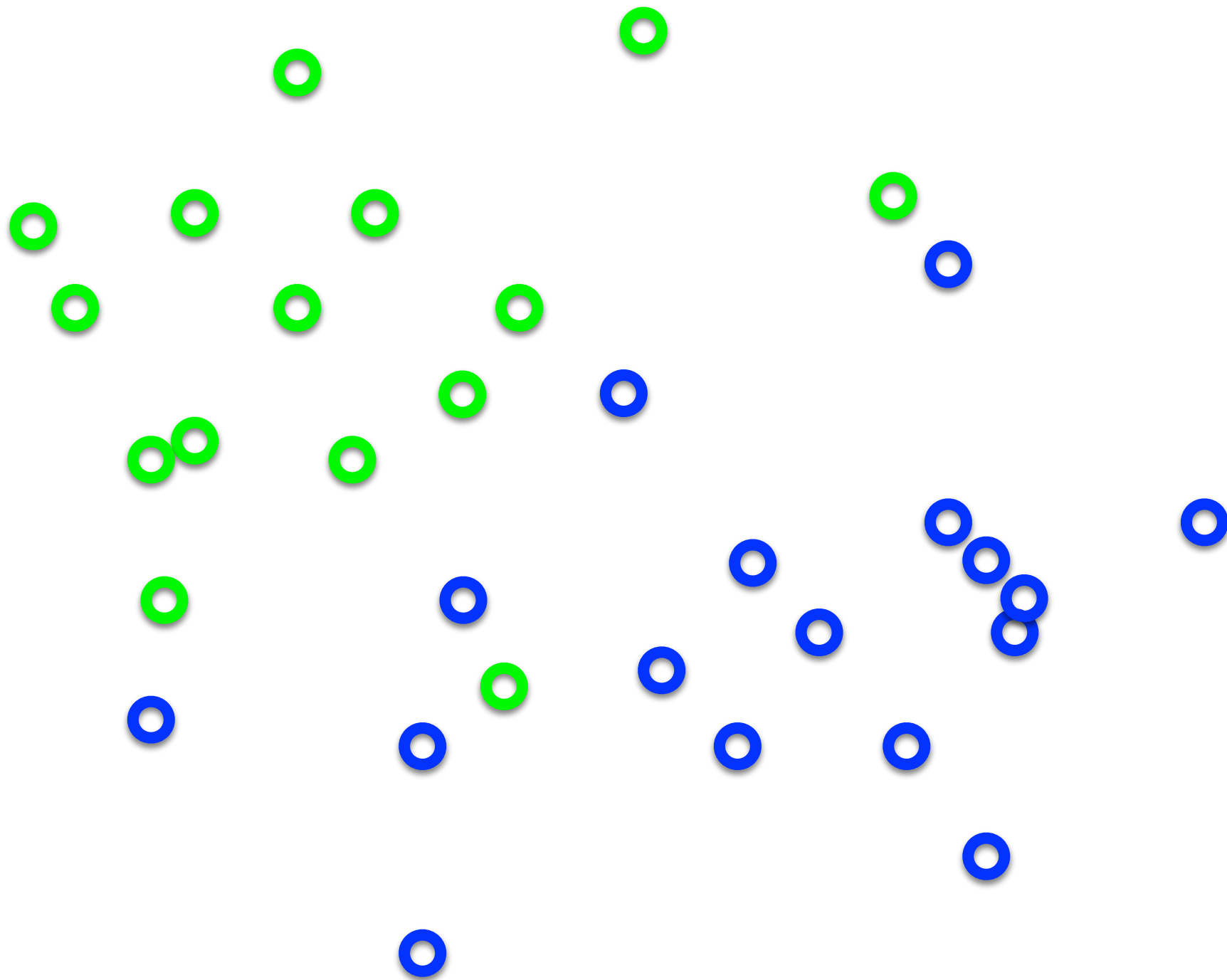


Parametric

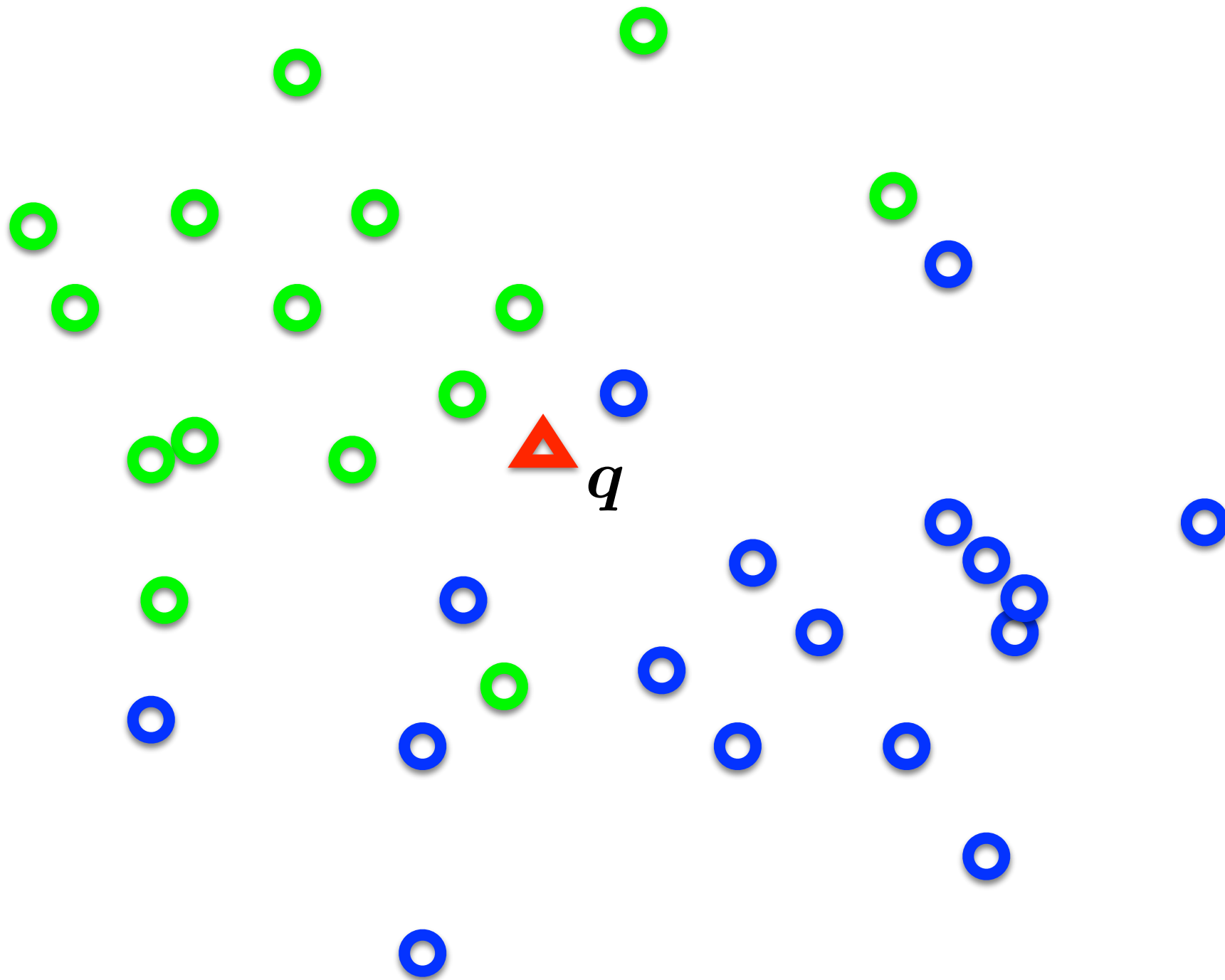
Discriminative



Distribution of data from two classes

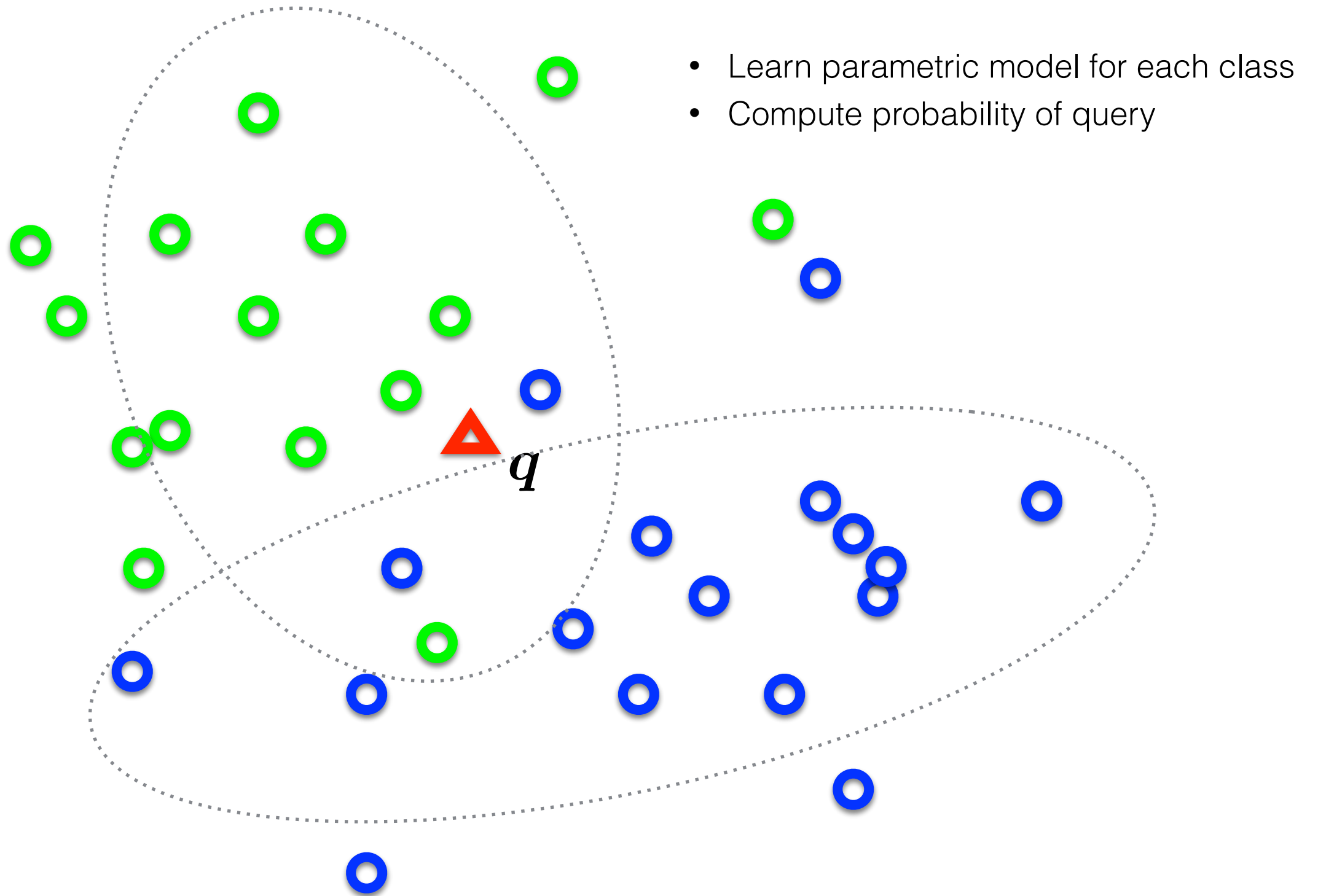


Distribution of data from two classes



Which class does q belong too?

Distribution of data from two classes



Let's consider the probability of q ...

This is called the posterior.

the probability of a class z given the observed features X

$$p(z|X)$$

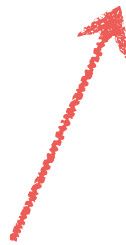
Recall: Bayes' Rule

$$\underset{\text{posterior}}{p(x|y)} = \frac{\overset{\text{likelihood}}{p(y|x)} \overset{\text{prior}}{p(x)}}{\underset{\text{evidence}}{p(y)}}$$

This is called the posterior.

the probability of a class z given the observed features X

$$p(z|X)$$

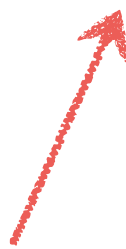


For classification, z is a
discrete random variable
(e.g., car, person, building)

This is called the posterior.

the probability of a class z given the observed features X

$$p(z|X)$$



For classification, z is a discrete random variable (e.g., car, person, building)



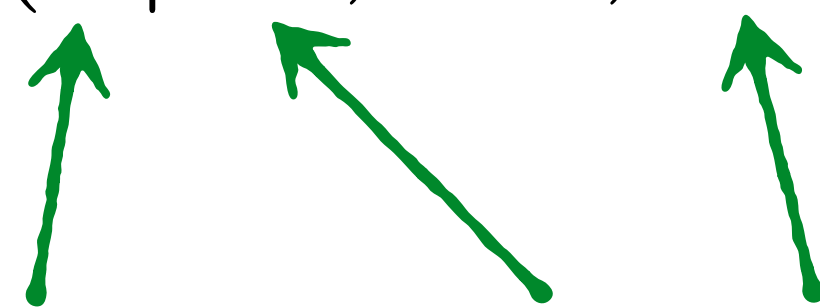
X is a **set** of observed feature (e.g., features from a single image)

(it's a function that returns a single probability value)

This is called the posterior.

the probability of a class z given the observed features X

$$p(z | x_1, \dots, x_N)$$



For classification, z is a
discrete random variable
(e.g., car, person, building)

Each x is an observed feature
(e.g., visual words)

(it's a function that returns a single probability value)

The **naive Bayes' classifier** is solving this optimization

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

MAP (maximum a posteriori) estimate

The naive Bayes' classifier is solving this optimization

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

MAP (maximum a posteriori) estimate

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} \frac{p(\mathbf{X} | z) p(z)}{p(\mathbf{X})}$$

Bayes' Rule

The naive Bayes' classifier is solving this optimization

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z|\mathbf{X})$$

MAP (maximum a posteriori) estimate

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} \frac{p(\mathbf{X}|z)p(z)}{p(\mathbf{X})}$$

Bayes' Rule

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(\mathbf{X}|z)p(z)$$

Remove constants

To optimize this...we need to compute this



How should we compute the likelihood? Make a naive assumption!

A naive Bayes' classifier assumes all features are
conditionally independent

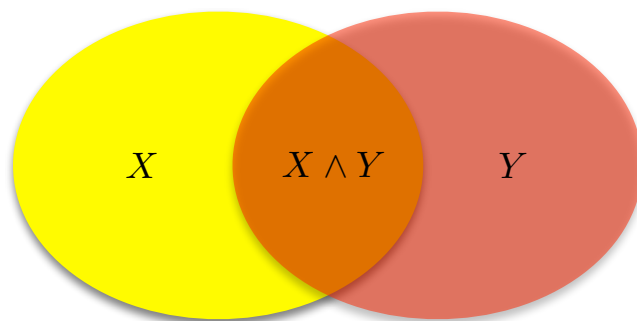
SO...

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(\mathbf{X}|z)p(z)$$

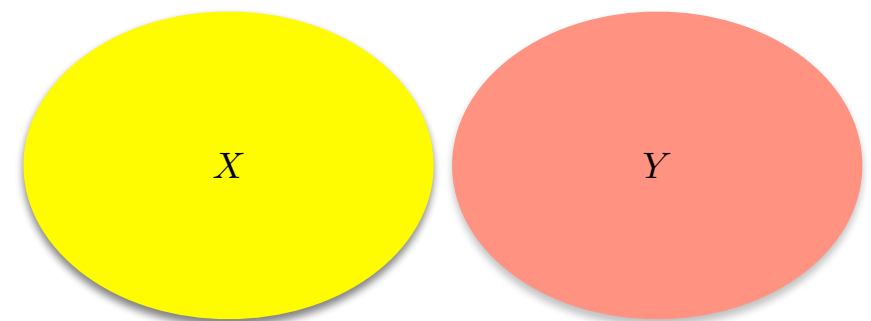
$$\begin{aligned} p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}) &= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2, \dots, \mathbf{x}_N | \mathbf{z}) \\ &= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) p(\mathbf{x}_3, \dots, \mathbf{x}_N | \mathbf{z}) \\ &= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) \cdots p(\mathbf{x}_N | \mathbf{z}) \end{aligned}$$

Each features z only depends on class x

Recall:



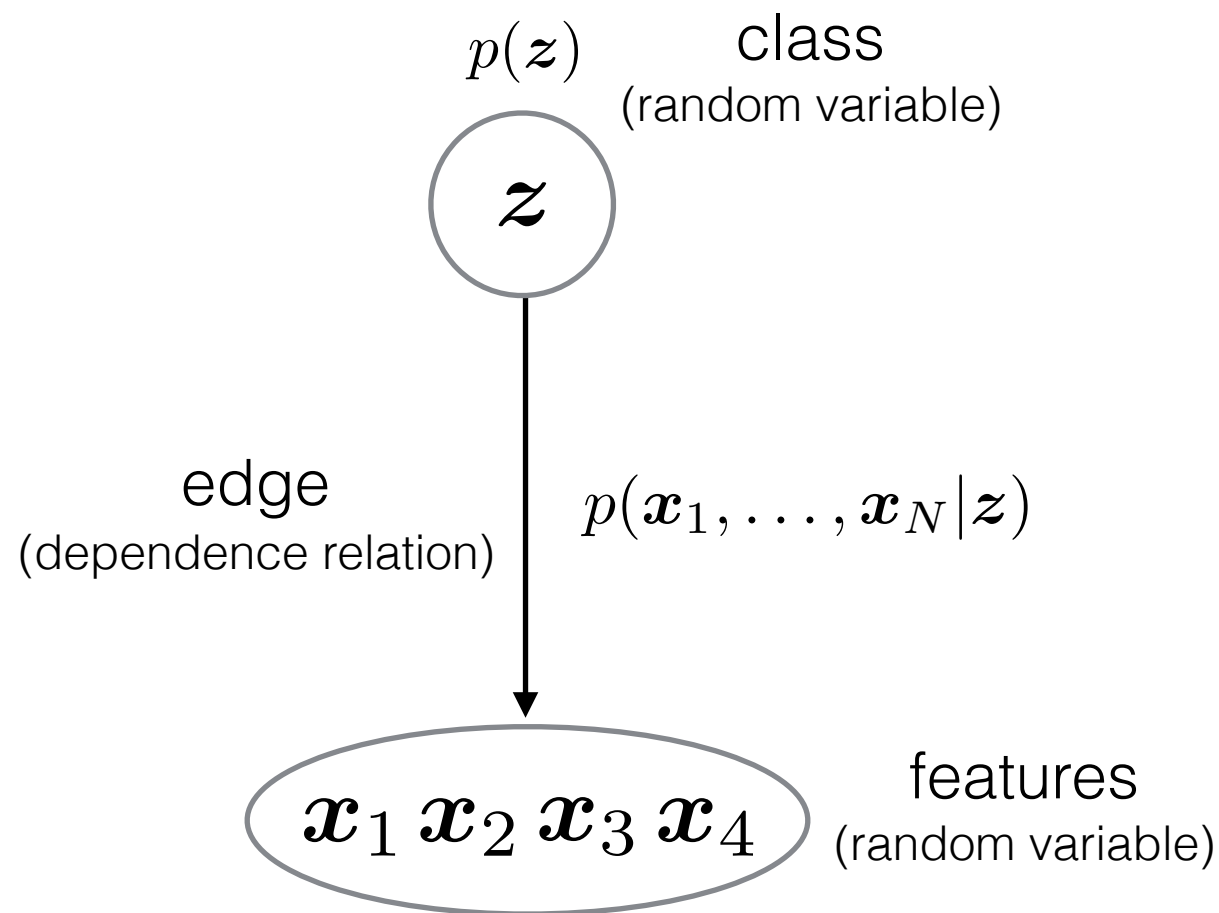
$$p(x, y) = p(x|y)p(y)$$



$$p(x, y) = p(x)p(y)$$

$$\begin{aligned}
p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}) &= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2, \dots, \mathbf{x}_N | \mathbf{z}) \\
&= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) p(\mathbf{x}_3, \dots, \mathbf{x}_N | \mathbf{z}) \\
&= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) \cdots p(\mathbf{x}_N | \mathbf{z})
\end{aligned}$$

You can visualize the distribution using a graphical model:

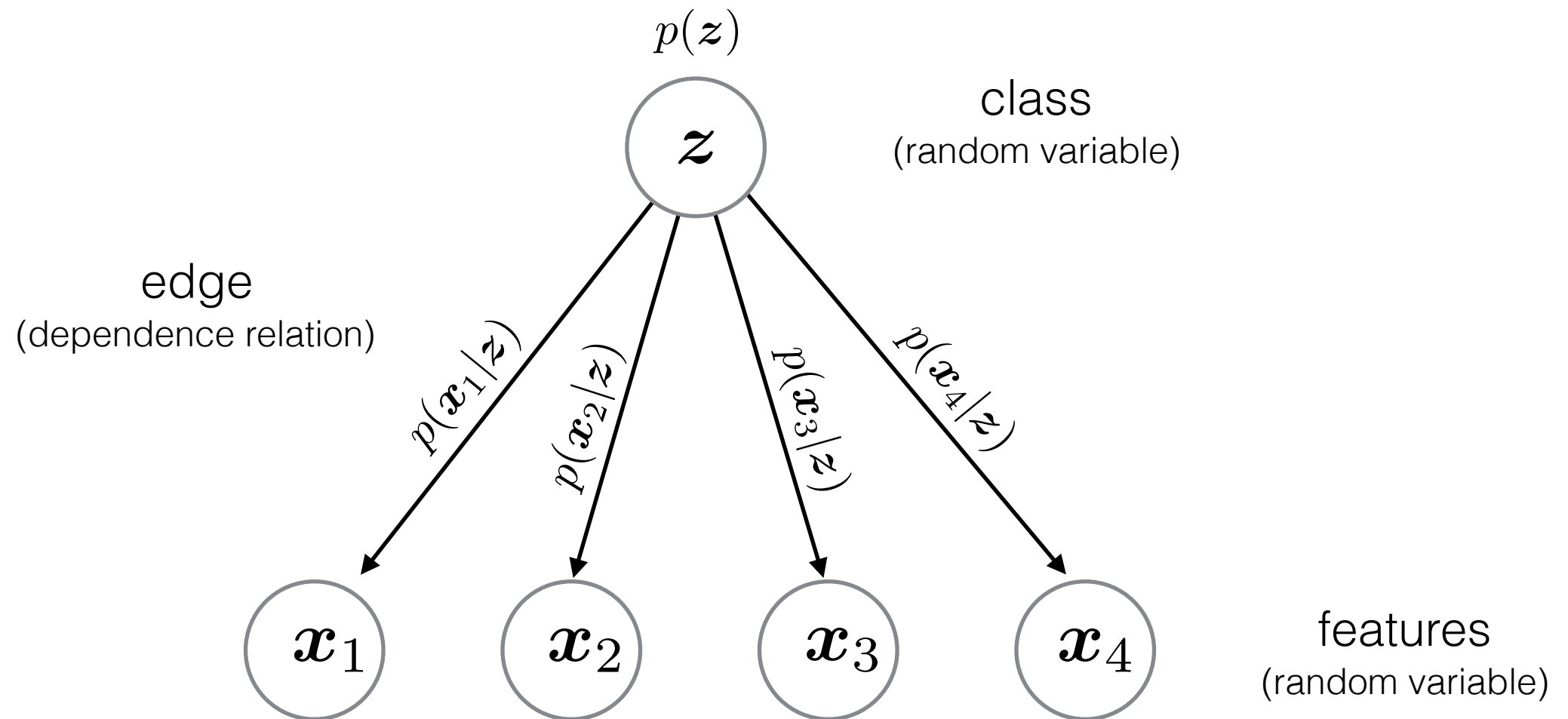


Instead of this. Do this ...

Naive Bayes assumption: If you assume conditional independence of the observations...

$$\begin{aligned}
p(\mathbf{x}_1, \dots, \mathbf{x}_N | \mathbf{z}) &= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2, \dots, \mathbf{x}_N | \mathbf{z}) \\
&= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) p(\mathbf{x}_3, \dots, \mathbf{x}_N | \mathbf{z}) \\
&= p(\mathbf{x}_1 | \mathbf{z}) p(\mathbf{x}_2 | \mathbf{z}) \cdots p(\mathbf{x}_N | \mathbf{z})
\end{aligned}$$

You can visualize the distribution using a graphical model:



To compute the MAP estimate

In order to solve this:

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

Given (1) a set of known parameters

$$p(\mathbf{z}) \quad p(\mathbf{x} | \mathbf{z})$$

(2) observations

$$\{x_1, x_2, \dots, x_N\}$$

To compute the MAP estimate

In order to solve this:

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

Given (1) a set of known parameters

$$p(\mathbf{z}) \quad p(\mathbf{x} | \mathbf{z})$$

(2) observations

$$\{x_1, x_2, \dots, x_N\}$$

Compute which z has the largest probability

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z) \prod_n p(x_n | z)$$

To compute the MAP estimate

In order to solve this:

$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z | \mathbf{X})$$

Given (1) a set of known parameters

$$p(\mathbf{z}) \quad p(\mathbf{x} | \mathbf{z})$$

(2) observations

$$\{x_1, x_2, \dots, x_N\}$$

Compute which z has the largest probability

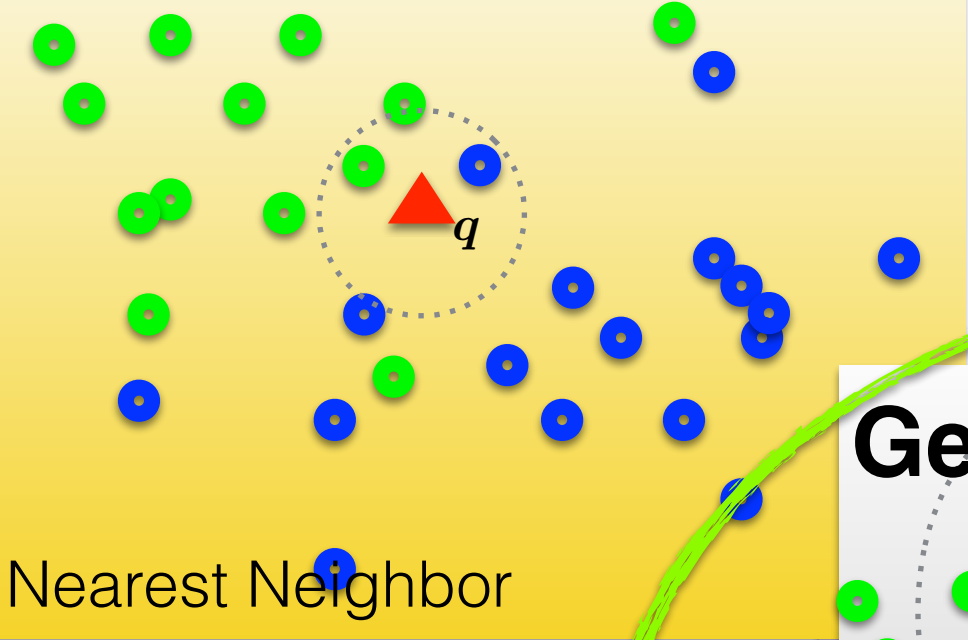
$$\hat{z} = \arg \max_{z \in \mathcal{Z}} p(z) \prod_n p(x_n | z)$$

This process is called '**inference**'

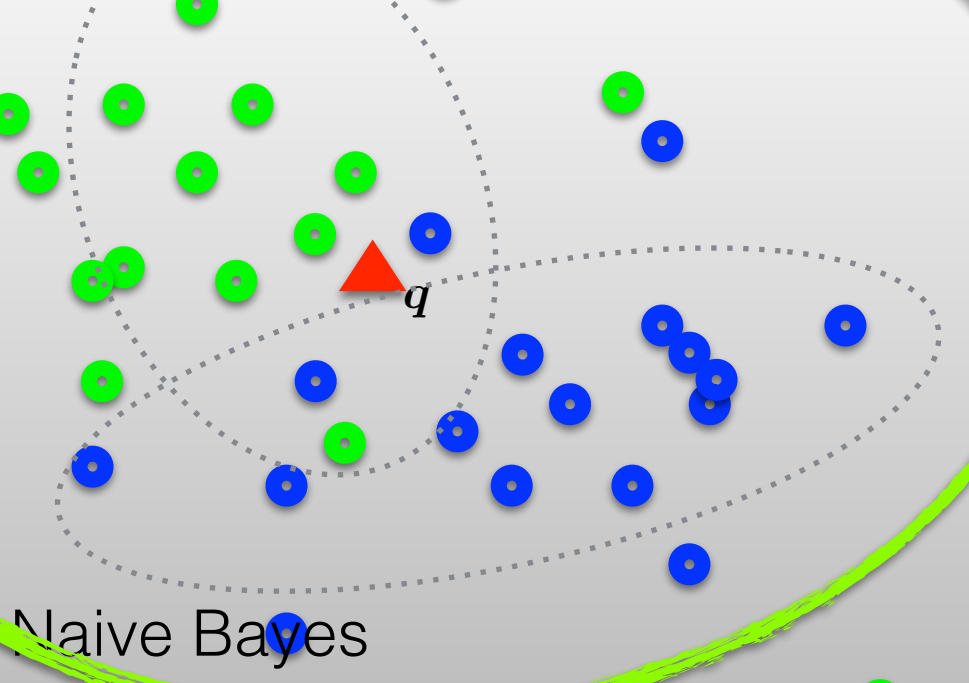
(you can learn more about 'learning' in machine learning)

(see concrete example at the end of slides)

Non-parametric

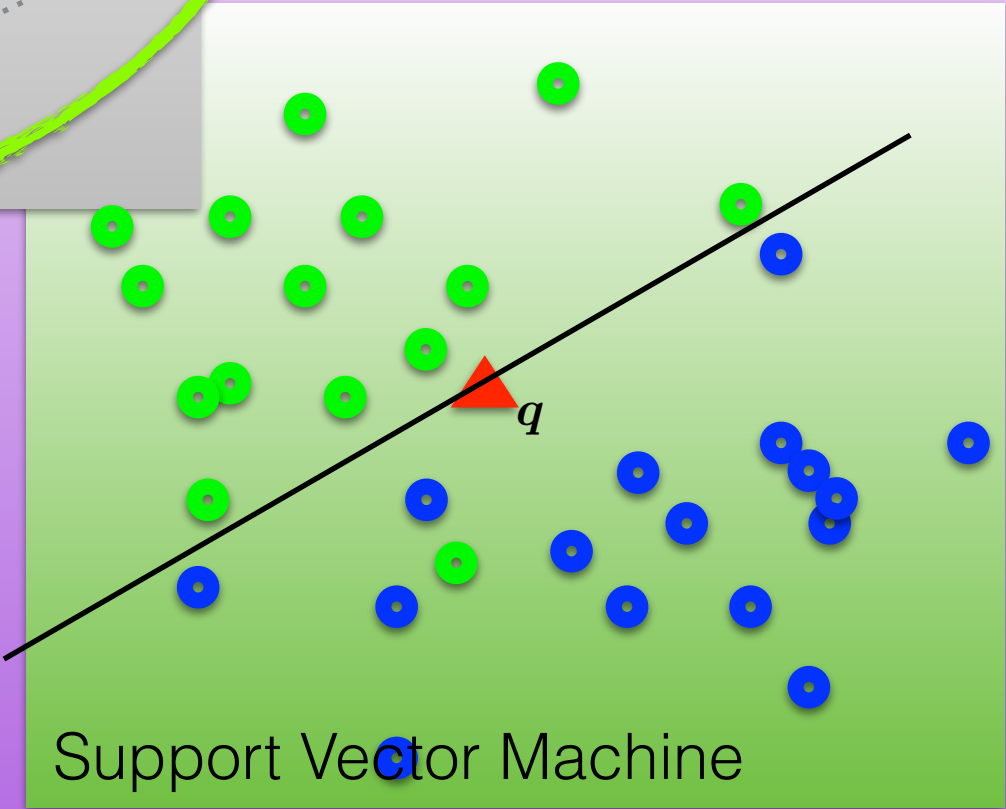


Generative



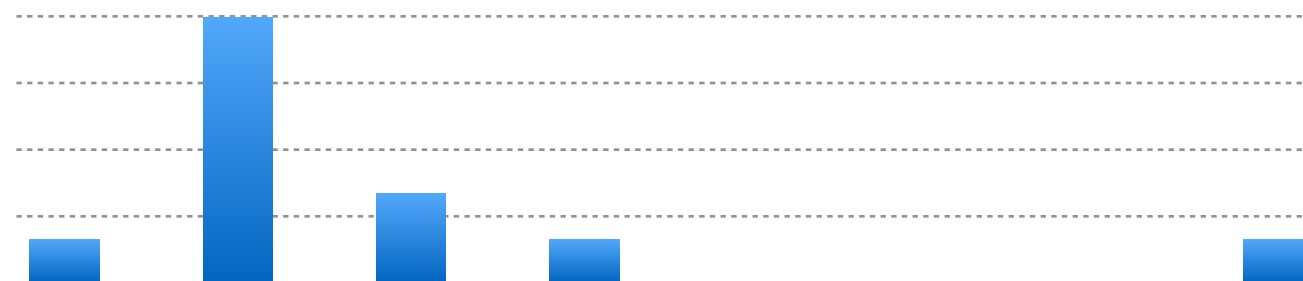
Parametric

Discriminative



Concrete Example:

Decomposing the likelihood
for the bag of words model

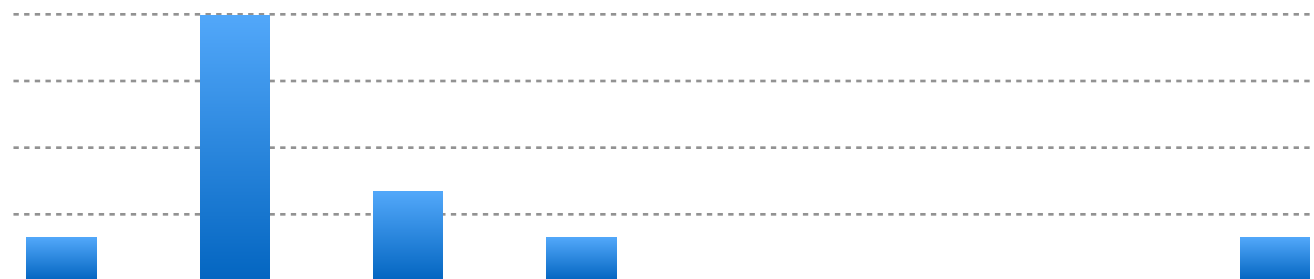


count	1	6	2	1	0	0	0	1
word	Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor
p(x z)	0.09	0.55	0.18	0.09	0.0	0.0	0.0	0.09

z= 'grand challenge'

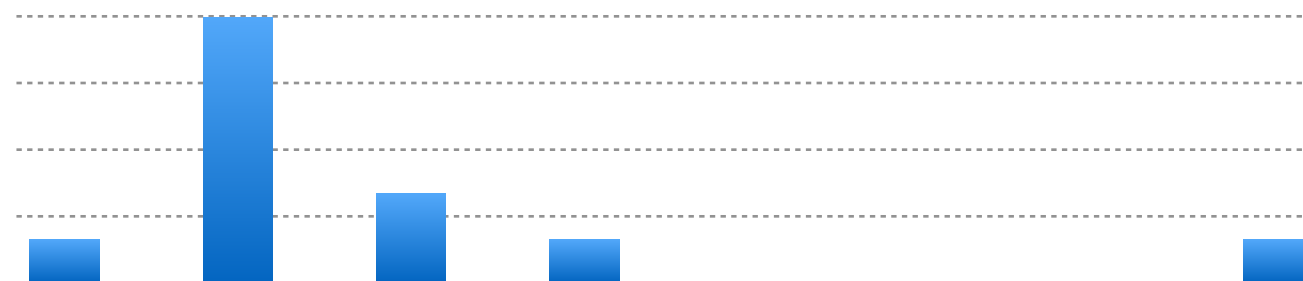
What is the likelihood of this document?

$$p(X|z)$$



count	1	6	2	1	0	0	0	1
word	Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor
p(x z)	0.09	0.55	0.18	0.09	0.0	0.0	0.0	0.09

$$\begin{aligned}
 p(X|z) &= \prod_v p(x_v|z)^{c(w_v)} \\
 &= (0.09)^1 (0.55)^6 \dots (0.09)^1
 \end{aligned}$$



count	1	6	2	1	0	0	0	1
word	Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor
p(x z)	0.09	0.55	0.18	0.09	0.0	0.0	0.0	0.09

$$p(X|z) = \prod_v p(x_v|z)^{c(w_v)}$$

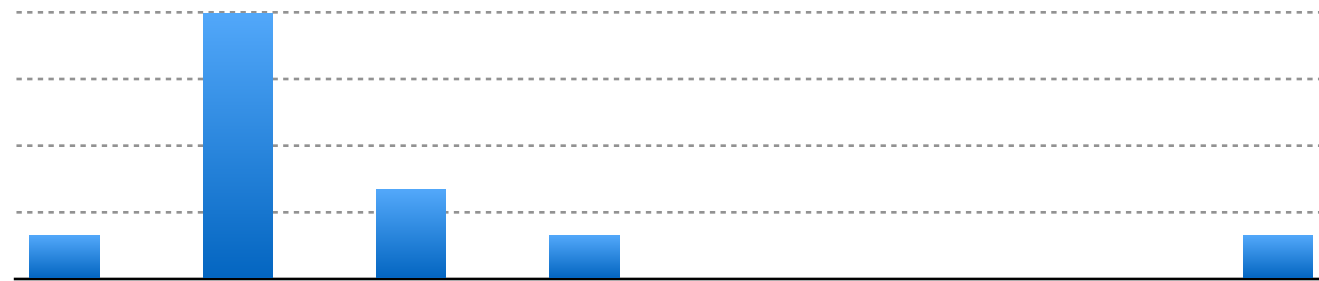
$$= (0.09)^1 (0.55)^6 \dots (0.09)^1$$

Numbers get really small so use log probabilities

$$\log p(X|z = \text{'grandchallenge'}) = -2.42 - 3.68 - 3.43 - 2.42 - 0.07 - 0.07 - 0.07 - 2.42 = -14.58$$

* typically add pseudo-counts (0.001)

** this is an example for computing the likelihood, need to multiply times **prior** to get posterior

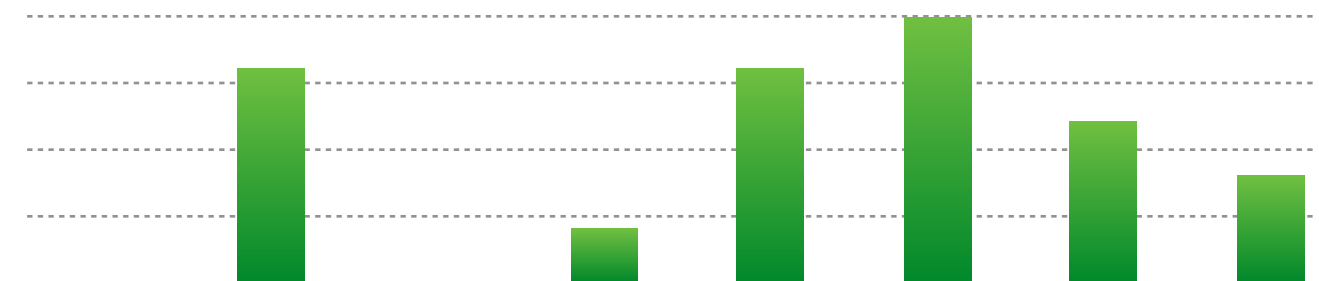


count	1	6	2	1	0	0	0	1
word	Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor
p(x z)	0.09	0.55	0.18	0.09	0.0	0.0	0.0	0.09

$$\log p(X|z=\text{grand challenge}) = - \mathbf{14.58}$$

z= 'grand challenge'

$$\log p(X|z=\text{bio inspired}) = - 37.48$$



count	0	4	0	1	4	5	3	2
word	Tartan	robot	CHIMP	CMU	bio	soft	ankle	sensor
p(x z)	0.0	0.21	0.0	0.05	0.21	0.26	0.16	0.11

$$\log p(X|z=\text{grand challenge}) = - 94.06$$

z= 'bio-inspired'

$$\log p(X|z=\text{bio inspired}) = - \mathbf{32.41}$$

* typically add pseudo-counts (0.001)

** this is an example for computing the likelihood, need to multiply times prior to get posterior