



Naive Bayes, MLE, MAP

Junxian He
Mar 10, 2026

Recap: Generative Models

Recap: Generative Models



X

$p(y | x)$

Discriminative

Cat

Y

Recap: Generative Models

$P(y)$ $P(x|y)$

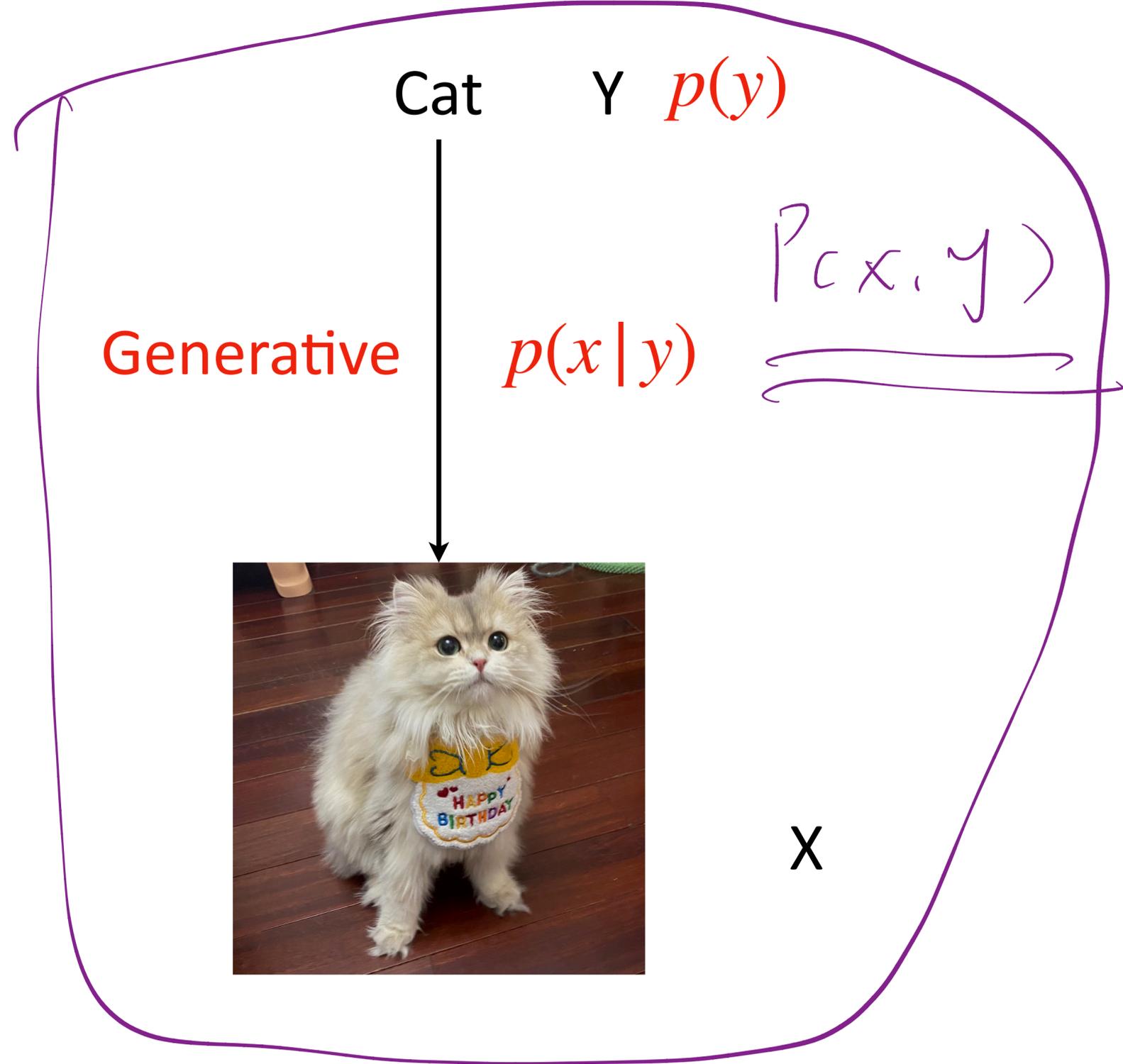


X

$p(y|x)$

Discriminative

Cat Y



Recap: Generative Models

posterior dist

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

→ prior

← $p(y|x)$

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

If our goal is to predict y , the distribution is often written as:

$p(y|x)$

→ prop to

$$p(y|x) \propto p(x|y)p(y)$$
$$\arg \max_y p(y|x) = \arg \max_y \frac{p(x|y)p(y)}{p(x)}$$
$$= \arg \max_y p(x|y)p(y)$$

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

How to represent the text?

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

How to represent the text?

if an email contains the j -th word of the dictionary, then we will set $x_j = 1$; otherwise, we let $x_j = 0$

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

How to represent the text?

if an email contains the j -th word of the dictionary, then we will set $x_j = 1$; otherwise, we let $x_j = 0$

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

How to represent the text?

if an email contains the j -th word of the dictionary, then we will set $x_j = 1$; otherwise, we let $x_j = 0$

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Dimension is the size of the dictionary

Recap: Naive Bayes

Binary classification: $y \in \{0,1\}$, x is discrete

Consider an email spam detection task, to predict whether the email is spam or not

How to represent the text?

if an email contains the j -th word of the dictionary, then we will set $x_j = 1$; otherwise, we let $x_j = 0$

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array} \quad \begin{array}{l} \text{vocabulary} \\ \\ \\ \text{Dimension is the size of the dictionary} \end{array}$$

Email Spam Classification

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Email Spam Classification

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

y spam or not x

Suppose the dictionary has 50000 words,
how many possible x?

$P(x|y)$ discrete

Email Spam Classification

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Suppose the dictionary has 50000 words,
how many possible x ?

2^{50000}
 $P(x|y)$
 $x = \dots$

Naive Bayes assumption: x_i 's are conditionally independent given y

Email Spam Classification

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Suppose the dictionary has 50000 words,
how many possible x ?

Naive Bayes assumption: x_i 's are conditionally independent given y

$$\text{For any } i \text{ and } j, \underline{p(x_i | y)} = \underline{p(x_i | y, x_j)}$$

Email Spam Classification

Email Spam Classification

$$\begin{aligned} & p(x_1, \dots, x_{50000} | y) \\ &= p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_1, x_2) \cdots p(x_{50000} | y, x_1, \dots, x_{49999}) \\ &= p(x_1 | y) p(x_2 | y) p(x_3 | y) \cdots p(x_{50000} | y) \\ &= \prod_{j=1}^d p(x_j | y) \end{aligned}$$

chain rule → autoregressive

conditionally independent

same

Email Spam Classification

$$p(x_1, \dots, x_{50000} | y)$$

Autoregressive

$$= p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_1, x_2) \cdots p(x_{50000} | y, x_1, \dots, x_{49999})$$

$$= p(x_1 | y) p(x_2 | y) p(x_3 | y) \cdots p(x_{50000} | y)$$

$$= \prod_{j=1}^d p(x_j | y)$$

Email Spam Classification

$$p(x_1, \dots, x_{50000} | y)$$

Autoregressive

$$= p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_1, x_2) \cdots p(x_{50000} | y, x_1, \dots, x_{49999})$$

$$= p(x_1 | y) p(x_2 | y) p(x_3 | y) \cdots p(x_{50000} | y)$$

$$= \prod_{j=1}^d p(x_j | y)$$

Parameters

$$\phi_{j|y=1} = p(x_j = 1 | y = 1), \quad \phi_{j|y=0} = p(x_j = 1 | y = 0), \quad \phi_y = p(y = 1)$$

Email Spam Classification

$$\begin{aligned} p(x_1, \dots, x_{50000} | y) & \quad \text{Autoregressive} \\ & = p(x_1 | y) p(x_2 | y, x_1) p(x_3 | y, x_1, x_2) \cdots p(x_{50000} | y, x_1, \dots, x_{49999}) \\ & = p(x_1 | y) p(x_2 | y) p(x_3 | y) \cdots p(x_{50000} | y) \\ & = \prod_{j=1}^d p(x_j | y) \end{aligned}$$

Parameters

$$\phi_{j|y=1} = p(x_j = 1 | y = 1), \quad \phi_{j|y=0} = p(x_j = 1 | y = 0), \quad \phi_y = p(y = 1)$$

50000 x 2 + 1 parameters (dict size is 50000) $\ll 2^{50000}$

Maximum Likelihood Estimation

Maximum Likelihood Estimation

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)})$$

Maximum Likelihood Estimation

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)})$$

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

$$\phi_y = \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}}{n}$$

Maximum Likelihood Estimation

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^n p(x^{(i)}, y^{(i)})$$

$$\begin{aligned} \phi_{j|y=1} &= \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}} \\ \phi_{j|y=0} &= \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}} \\ \phi_y &= \frac{\sum_{i=1}^n 1\{y^{(i)} = 1\}}{n} \end{aligned}$$

Count the occurrence of x_j in spam/
non-spam emails and normalize

$\boxed{P(y=1)}$ \rightarrow # emails

Prediction

Prediction

$$\begin{aligned} p(y = 1|x) &= \frac{p(x|y = 1)p(y = 1)}{p(x)} \\ &= \frac{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1) + \left(\prod_{j=1}^d p(x_j|y = 0)\right) p(y = 0)} \end{aligned}$$

Prediction

$$\begin{aligned} p(y = 1|x) &= \frac{p(x|y = 1)p(y = 1)}{p(x)} \\ &= \frac{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1) + \left(\prod_{j=1}^d p(x_j|y = 0)\right) p(y = 0)} \end{aligned}$$

Naive Classifier

$P(y=1|x)$

$P(y=0|x)$

Laplace Smoothing

Laplace Smoothing

What if we never see the word “learning” in training data but “learning” exists in the test data?

Laplace Smoothing

What if we never see the word “learning” in training data but “learning” exists in the test data?

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}$$
$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$


$x_j =$ “learning”

Laplace Smoothing

What if we never see the word “learning” in training data but “learning” exists in the test data?

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

Suppose the index in the dictionary for “learning” is q

$$p(x_q = 1 | y = 1) = 0$$

$$p(x_q = 1 | y = 0) = 0$$



Laplace Smoothing

What if we never see the word “learning” in training data but “learning” exists in the test data?

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

Suppose the index in the dictionary for “learning” is q

$$p(x_q = 1 | y = 1) = 0$$

$$p(x_q = 1 | y = 0) = 0$$

$$p(y = 1 | x) = \frac{p(x|y = 1)p(y = 1)}{p(x)}$$

$$= \frac{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1) + \left(\prod_{j=1}^d p(x_j|y = 0)\right) p(y = 0)}$$

Laplace Smoothing

What if we never see the word “learning” in training data but “learning” exists in the test data?

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^n 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

Suppose the index in the dictionary for “learning” is q

$$p(x_q = 1 | y = 1) = 0$$

$$p(x_q = 1 | y = 0) = 0$$

$$p(y = 1 | x) = \frac{p(x|y = 1)p(y = 1)}{p(x)}$$

$$= \frac{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1)}{\left(\prod_{j=1}^d p(x_j|y = 1)\right) p(y = 1) + \left(\prod_{j=1}^d p(x_j|y = 0)\right) p(y = 0)}$$

$$= \frac{0}{0}$$



Laplace Smoothing

Laplace Smoothing

Take the problem of estimating the mean of a multinomial random variable z taking values in $\{1, \dots, k\}$. Given the independent observations $\{z^{(1)}, \dots, z^{(n)}\}$



Laplace Smoothing

Take the problem of estimating the mean of a multinomial random variable z taking values in $\{1, \dots, k\}$. Given the independent observations $\{z^{(1)}, \dots, z^{(n)}\}$

$$\phi_j = p(z = j)$$

$$\phi_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}}{n}$$

Laplace Smoothing

Take the problem of estimating the mean of a multinomial random variable z taking values in $\{1, \dots, k\}$. Given the independent observations $\{z^{(1)}, \dots, z^{(n)}\}$

$$\phi_j = p(z = j) \qquad \phi_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}}{n}$$

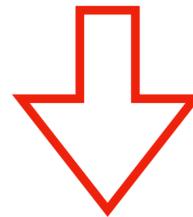
A handwritten diagram illustrating the Laplace smoothing formula. A large purple oval encloses the formula $\phi_j = \frac{1 + \sum_{i=1}^n 1\{z^{(i)} = j\}}{k + n}$. A red arrow points down from the second formula above to the '1' in the numerator of this formula. A purple arrow points from the '0' in the denominator of the second formula above to the 'k' in the denominator of this formula. The '1' in the numerator is circled in purple, and the 'k' in the denominator is also circled in purple. The word 'Laplace' is written in purple above the '1'.

$$\phi_j = \frac{1 + \sum_{i=1}^n 1\{z^{(i)} = j\}}{k + n}$$

Laplace Smoothing

Take the problem of estimating the mean of a multinomial random variable z taking values in $\{1, \dots, k\}$. Given the independent observations $\{z^{(1)}, \dots, z^{(n)}\}$

$$\phi_j = p(z = j) \qquad \phi_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}}{n}$$



$$\phi_j = \frac{1 + \sum_{i=1}^n 1\{z^{(i)} = j\}}{k + n}$$

Why adding k to the denominator?

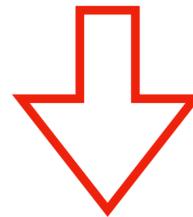
$$\sum_{j=1}^k \phi_j = 1 \qquad \frac{n+k}{n+k}$$

Laplace Smoothing

Take the problem of estimating the mean of a multinomial random variable z taking values in $\{1, \dots, k\}$. Given the independent observations $\{z^{(1)}, \dots, z^{(n)}\}$

$$\phi_j = p(z = j)$$

$$\phi_j = \frac{\sum_{i=1}^n 1\{z^{(i)} = j\}}{n}$$



$$\phi_j = \frac{1 + \sum_{i=1}^n 1\{z^{(i)} = j\}}{k + n}$$

Why adding k to the denominator?

In the email spam classification case:

$$\phi_{j|y=1} = \frac{1 + \sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{2 + \sum_{i=1}^n 1\{y^{(i)} = 1\}}$$
$$\phi_{j|y=0} = \frac{1 + \sum_{i=1}^n 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{2 + \sum_{i=1}^n 1\{y^{(i)} = 0\}}$$

MSE

$$\frac{1}{2} \sum (\hat{y} - y)^2$$

why?

$$\sum \frac{1}{\sigma} | \hat{y} - y |$$

$$\sum \frac{1}{\sigma} (\hat{y} - y)^4$$

equivalent

MLE Gaussian

why MLE?

Parameter Estimation: MLE and MAP

$$E \left[\frac{\hat{y}}{y} - 1 \right]^2$$

Maximum Likelihood Estimation (MLE)

dataset distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

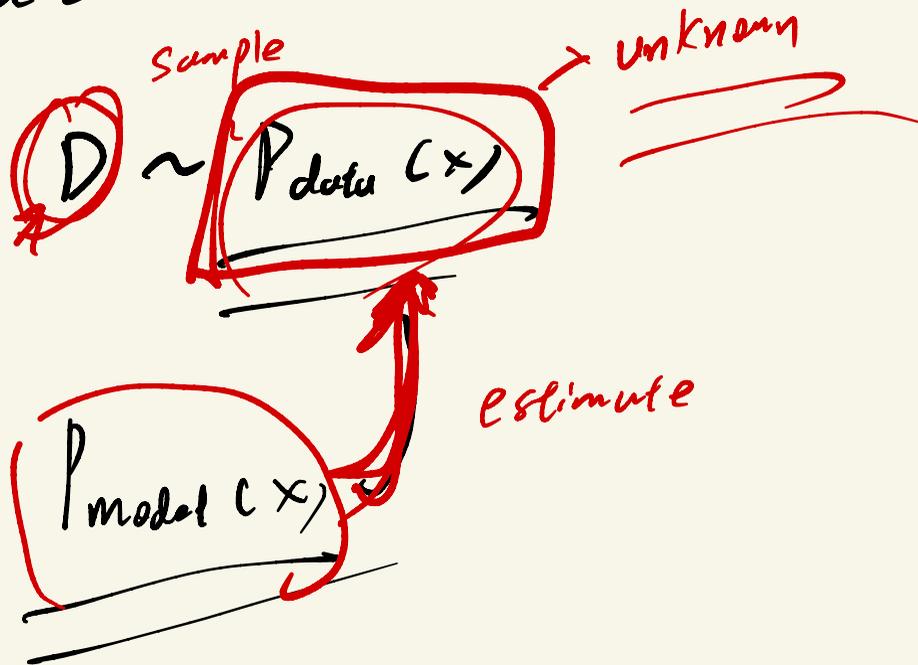
parameters

In practice:

$$\arg \max_{\theta} \frac{1}{n} \sum_i \log p_{model}(x^{(i)}; \theta)$$

$x^{(1)}, x^{(2)}, x^{(n)} \sim P_{data}(x)$

$D = \text{dataset}$



Maximum Likelihood Estimation (MLE)

Suppose $p_{data}(x)$ is the real data distribution, $p_{model}(x; \theta)$ is our model parameterized by θ

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

In practice:

$$\arg \max_{\theta} \frac{1}{n} \sum_i^n \log p_{model}(x^{(i)}; \theta)$$

Maximum Likelihood Estimation (MLE)

Suppose $p_{data}(x)$ is the real data distribution, $p_{model}(x; \theta)$ is our model parameterized by θ

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

In practice:

$$\arg \max_{\theta} \frac{1}{n} \sum_i \log p_{model}(x^{(i)}; \theta)$$

relationship

$x^{(i)}$ are i.i.d. (independent and identically distributed) samples from $p_{data}(x)$

Central Limit Theorem

$n \rightarrow \infty$

$$\frac{1}{n} \sum_i \log p_{model}(x^{(i)}; \theta) \rightarrow \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

Maximum Likelihood Estimation (MLE)

Suppose $p_{data}(x)$ is the real data distribution, $p_{model}(x; \theta)$ is our model parameterized by θ

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

In practice:

$$\arg \max_{\theta} \frac{1}{n} \sum_i \log p_{model}(x^{(i)}; \theta)$$

$x^{(i)}$ are i.i.d. (independent and identically distributed) samples from $p_{data}(x)$

Monte Carlo Estimation of Expectation

Maximum Likelihood Estimation (MLE)

Suppose $p_{data}(x)$ is the real data distribution, $p_{model}(x; \theta)$ is our model parameterized by θ

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

In practice:

$$\arg \max_{\theta} \frac{1}{n} \sum_i \log p_{model}(x^{(i)}; \theta)$$

θ

$x^{(i)}$ are i.i.d. (independent and identically distributed) samples from $p_{data}(x)$

Monte Carlo Estimation of Expectation

Why can we make this approximation?

unbiased estimation:

estimation: $\hat{\theta}$ θ^*

$$E(\hat{\theta}) = \theta^*$$

unbiased

$n \rightarrow \infty$

Monte Carlo Estimation of Expectation

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \mathbb{E}_{x \sim p(x)} f(x) \rightarrow \text{real unbiased}$$
$$\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \frac{\text{Var}(f(x))}{n}$$

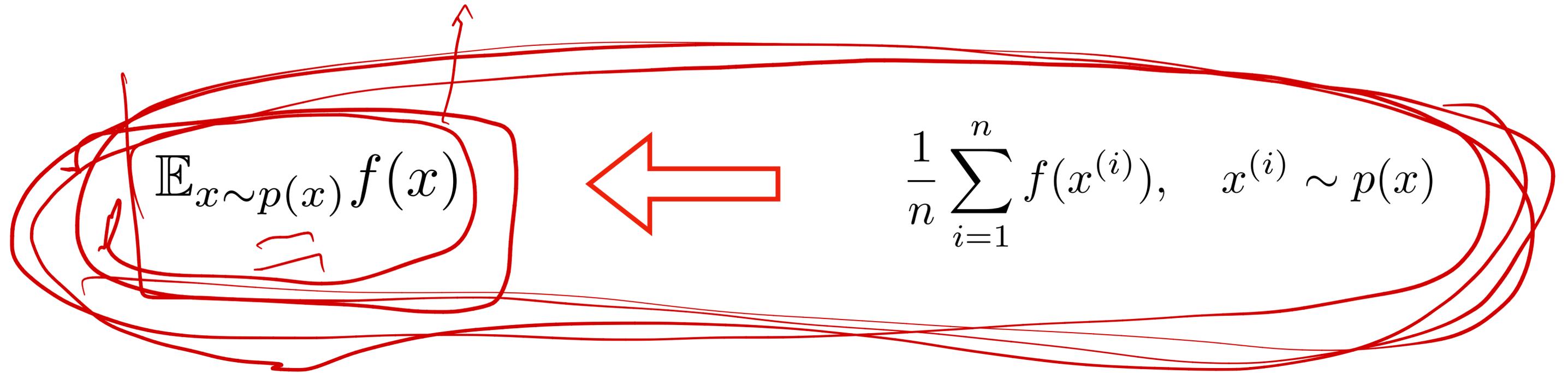
Monte Carlo Estimation of Expectation

$$\mathbb{E}_{x \sim p(x)} f(x)$$

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \mathbb{E}_{x \sim p(x)} f(x)$$

$$\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \frac{\text{Var}(f(x))}{n}$$

Monte Carlo Estimation of Expectation



$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \mathbb{E}_{x \sim p(x)} f(x)$$

$$\text{Var}\left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)})\right] = \frac{\text{Var}(f(x))}{n}$$

Monte Carlo Estimation of Expectation

$$\mathbb{E}_{x \sim p(x)} f(x)$$

$$\frac{1}{n} \sum_{i=1}^n f(x^{(i)}), \quad x^{(i)} \sim p(x)$$

In practice, n is often small, like 1 sample

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)}) \right] = \mathbb{E}_{x \sim p(x)} f(x)$$

$$\text{Var} \left[\frac{1}{n} \sum_{i=1}^n f(x^{(i)}) \right] = \frac{\text{Var}(f(x))}{n}$$

Sampling and Evaluation of Distributions

Sampling and Evaluation of Distributions

- Some distributions are easy to sample from but hard to compute the probability value (hard to evaluate)
- Monte Carlo estimation requires this kind of distribution

$P(x)$

hard to sample but easy to compute

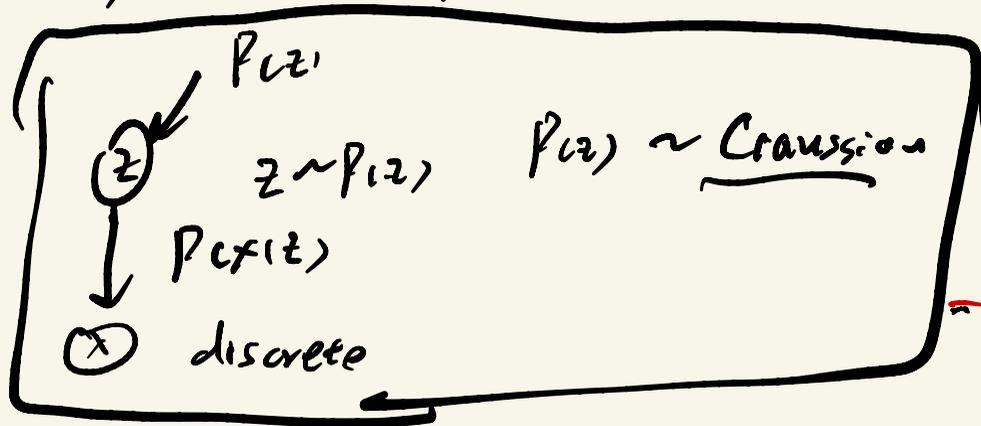
* continuous

$$P(x) = \frac{e^{x+1}}{\dots}$$

x_0

$$\int_z P(x_0, z) P(z) dz$$

easy to sample but hard to compute



$$P(x)$$

$x \sim P(x)$

$z_0 \sim P(z), x_0 \sim P(x|z_0)$

$$x_1 \int_z P(x_1, z) dz$$

Sampling and Evaluation of Distributions

- Some distributions are easy to sample from but hard to compute the probability value (hard to evaluate)
 - Monte Carlo estimation requires this kind of distribution
- Some distributions are easy to compute the probability value (easy to evaluate) but hard to sample from
 - How to sample from a distribution efficiently is a separate topic

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

What is the optimal p_{model} ?

$$p_{model} \longrightarrow p_{data}(x)$$

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

What is the optimal p_{model} ?

MLE is equivalent to

$$\arg \min_{\theta} D_{KL}(p_{data}(x) || p_{model}(x; \theta))$$

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

What is the optimal p_{model} ?

MLE is equivalent to

KL divergence

$$\arg \min_{\theta} D_{KL}(p_{data}(x) || p_{model}(x; \theta))$$

$D_{KL} \geq 0$ is a distance metric between two distributions, it is 0 when the two distributions are identical

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

What is the optimal p_{model} ?

MLE is equivalent to

$$\arg \min_{\theta} D_{KL}(p_{data}(x) || p_{model}(x; \theta))$$

$D_{KL} \geq 0$ is a distance metric between two distributions, it is 0 when the two distributions are identical

$$D_{KL}(p(x) || q(x)) = \mathbb{E}_{p(x)} \log \frac{p(x)}{q(x)}$$

MLE is Approximating the Real Distribution

$$\arg \max_{\theta} \mathbb{E}_{x \sim p_{data}(x)} \log p_{model}(x; \theta)$$

What is the optimal p_{model} ?

MLE is equivalent to

$$\arg \min_{\theta} D_{KL}(p_{data}(x) || p_{model}(x; \theta))$$

$D_{KL} \geq 0$ is a distance metric between two distributions, it is 0 when the two distributions are identical

$$D_{KL}(p(x) || q(x)) = \mathbb{E}_{p(x)} \log \frac{p(x)}{q(x)}$$

When data is all the data from the world, then MLE is learning a distribution for the world

$$KL(P \parallel q) = E_p \left(\log \frac{P}{q} \right) \quad \log P - \log q$$

$$\operatorname{argmin}_{\theta} KL [P_{\text{data}}(x) \parallel P_{\text{model}}(x; \theta)]$$

$$= \operatorname{argmin}_{\theta} \left[E_{P_{\text{data}}(x)} \left(\log P_{\text{data}}(x) \right) - \log P_{\text{model}}(x; \theta) \right]$$

$$= \operatorname{argmin}_{\theta} -E_{P_{\text{data}}(x)} \log P_{\text{model}}(x; \theta) = \operatorname{argmax}_{\theta} E_{P_{\text{data}}(x)} \log P_{\text{model}}(x; \theta)$$

MLE

KL divergence

$$KL(P \parallel Q) = \mathbb{E}_P \left[\log \frac{p}{q} \right] \geq 0$$

~~Jensen Inequality~~ \Rightarrow when $P(x) = Q(x)$

$$-\mathbb{E}_P \left(\log \frac{q}{p} \right) \geq -\log \left(\mathbb{E}_P \frac{q}{p} \right)$$

concave

$$\mathbb{E}_P \frac{q}{p} = \int_x P(x) \cdot \frac{q(x)}{P(x)}$$

$$\underline{q(x) = P(x) \text{ everywhere} = 1}$$

JSD

$$KL(P_{data} \parallel P_{model}) = E_{P_{data}(x)}$$

$$\log \frac{P_{data}}{P_{model}}$$

$x \sim P_{data}(x)$

$$KL(P_{model} \parallel P_{data}) = E_{P_{model}(x; \theta)}$$

$$\log \frac{P_{model}}{P_{data}}$$

depend on θ

$$x \sim P_{model}(x; \theta)$$

Biased/Unbiased Estimator

Biased/Unbiased Estimator

Suppose we want to estimate a true quantity θ^* , and our estimation is $\hat{\theta}$, then we define the bias of the estimation as:

Biased/Unbiased Estimator

Suppose we want to estimate a true quantity θ^* , and our estimation is $\hat{\theta}$, then we define the bias of the estimation as:

$$\text{bias} = \mathbb{E}(\hat{\theta}) - \theta^*$$

Biased/Unbiased Estimator

Suppose we want to estimate a true quantity θ^* , and our estimation is $\hat{\theta}$, then we define the bias of the estimation as:

$$\text{bias} = \mathbb{E}(\hat{\theta}) - \theta^*$$

When does the estimation converges to the true value when we have infinite data samples?

Biased/Unbiased Estimator

Suppose we want to estimate a true quantity θ^* , and our estimation is $\hat{\theta}$, then we define the bias of the estimation as:

$$bias = \mathbb{E}(\hat{\theta}) - \theta^*$$

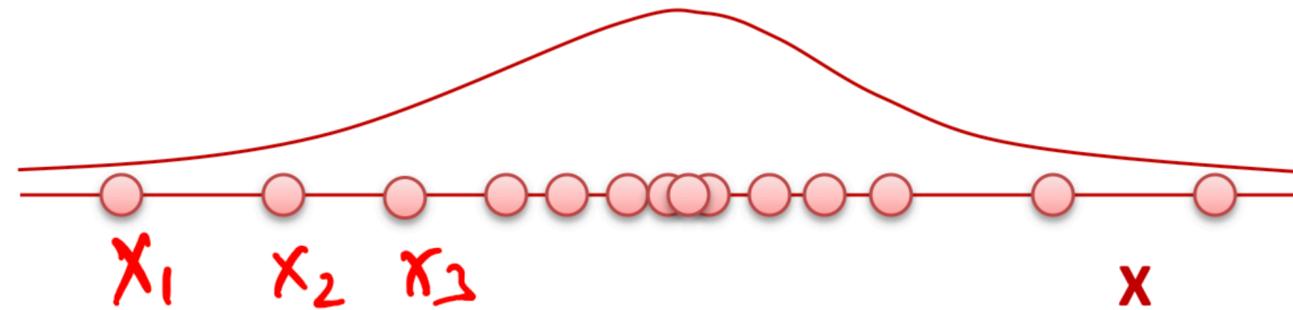
When does the estimation converges to the true value when we have infinite data samples?

$$bias \rightarrow 0,$$

$$Var(\hat{\theta}) \rightarrow 0$$

Learn Parameters from Data with MLE

Data, $D =$



Approximate the mean and variance of the data

Data are **i.i.d.:**

- **Independent** events
- **Identically distributed** according to Gaussian distribution

MLE for Gaussian Mean and Variance

$$E\left[\frac{1}{n} \sum_{i=1}^n x_i\right]$$

$$= \frac{\sum_{i=1}^n E(x_i)}{n} \quad \mathcal{M}$$

\mathcal{M} unbiased

$$\hat{\mu}_{MLE}$$

$$= \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}_{MLE}^2$$

$$= \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

$$E[(x - \mu)^2]$$



MLE for Gaussian Mean and Variance

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

Are the estimations biased?

MLE for Gaussian Mean and Variance

unbiased $\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$

MLE $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2$ *biased*

$E[\hat{\sigma}_{MLE}^2] = \sigma^2$

$= \frac{n-1}{n} \sigma^2$

Are the estimations biased?

Unbiased estimator: $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$

not MLE

Max A Posterior (MAP) Estimation

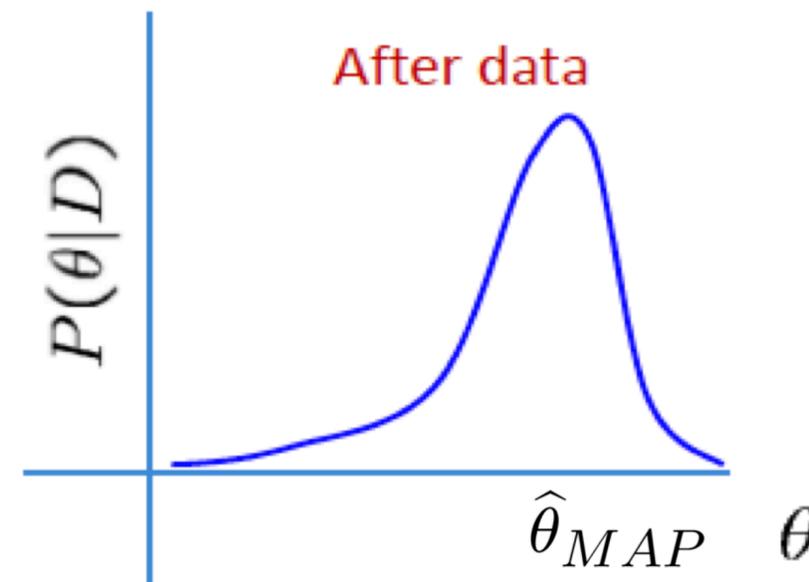
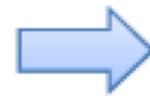
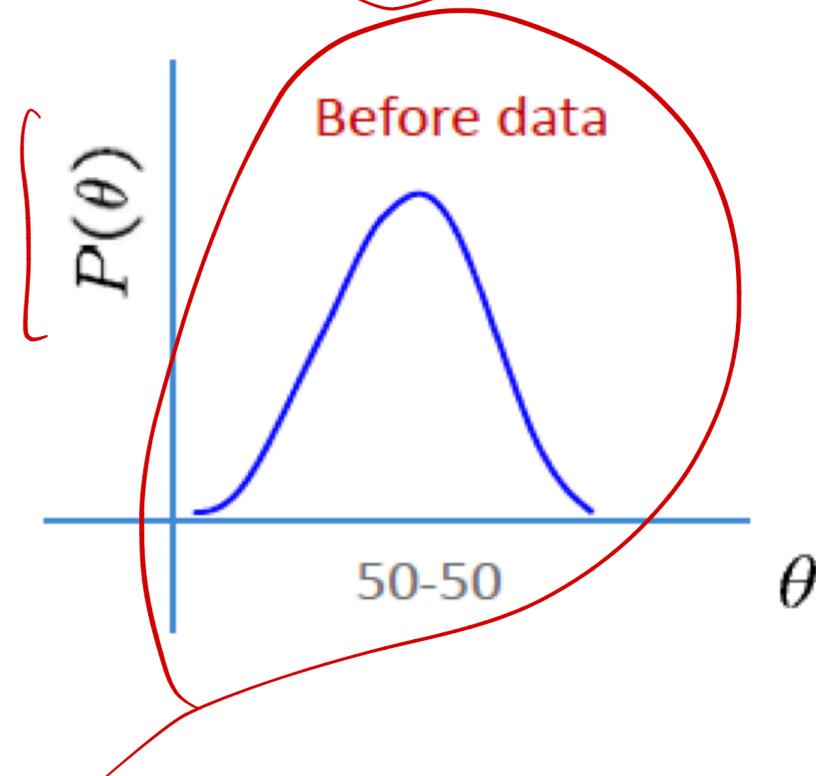
Max A Posterior (MAP) Estimation

Bring prior knowledge to the parameter, define the prior $P(\theta)$. The posterior distribution is $P(\theta | D)$. D is the training dataset

$$P(\theta) \sim N(0, 1)$$

Max A Posterior (MAP) Estimation

Bring prior knowledge to the parameter, define the prior $P(\theta)$. The posterior distribution is $P(\theta | D)$. D is the training dataset

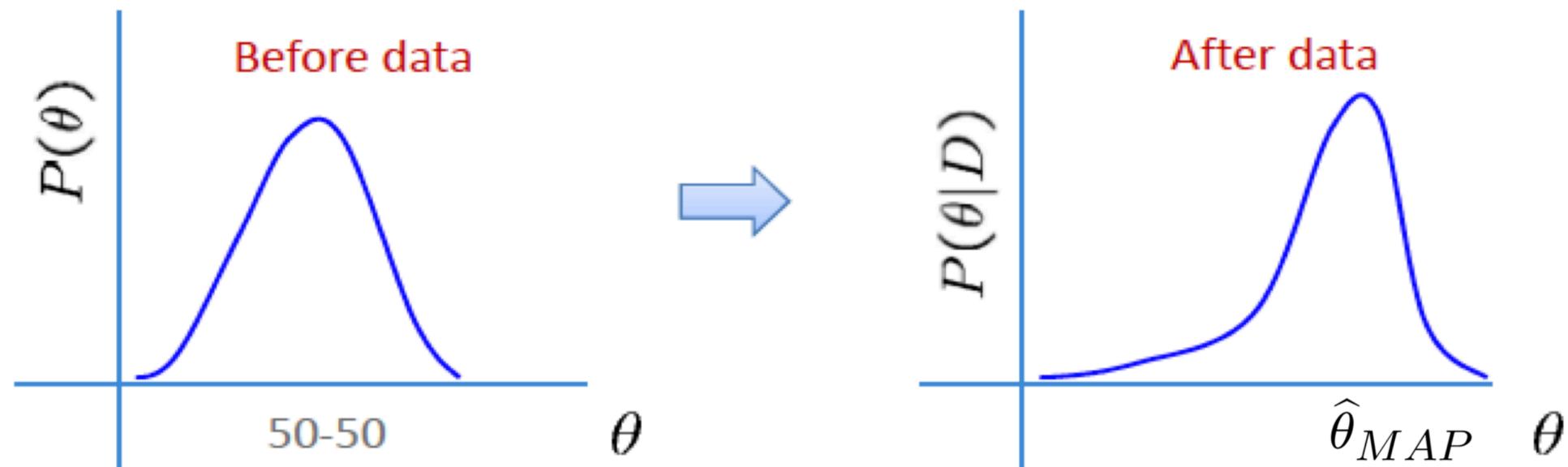


~~$P(\theta|D) \propto P(\theta)P(D|\theta)$~~
 $P(\theta)P(D|\theta)$

Handwritten blue annotations include a circle around $P(\theta)$, a larger circle around $P(D|\theta)$, and a blue arrow pointing downwards from the right side of the diagram.

Max A Posterior (MAP) Estimation

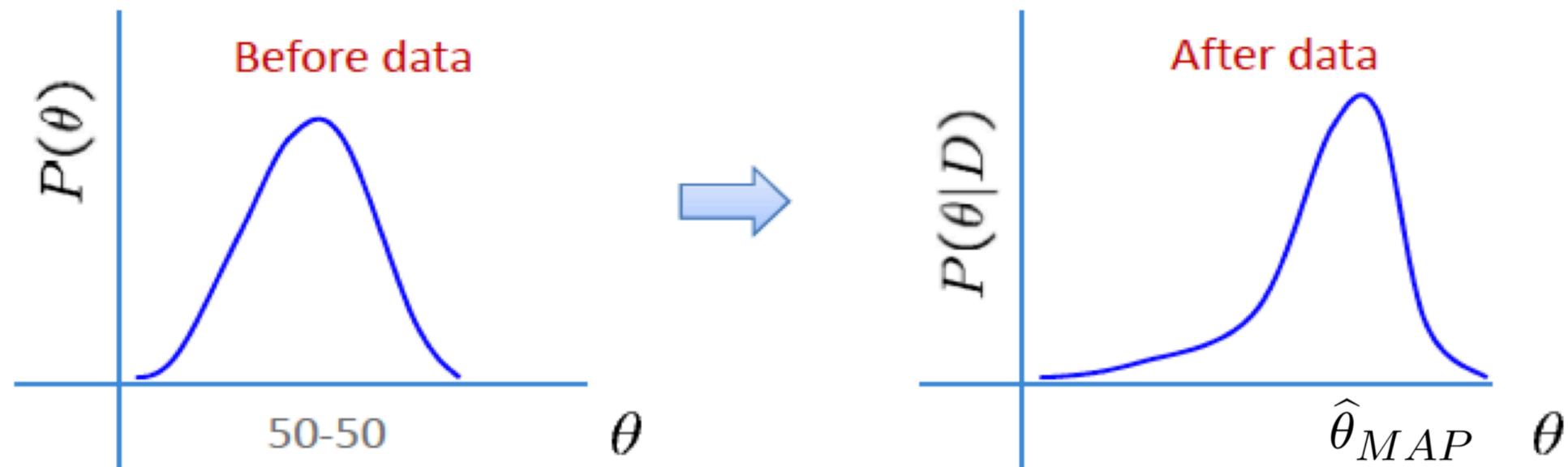
Bring prior knowledge to the parameter, define the prior $P(\theta)$. The posterior distribution is $P(\theta | D)$. D is the training dataset



$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} P(D | \theta) P(\theta)\end{aligned}$$

Max A Posterior (MAP) Estimation

Bring prior knowledge to the parameter, define the prior $P(\theta)$. The posterior distribution is $P(\theta | D)$. D is the training dataset

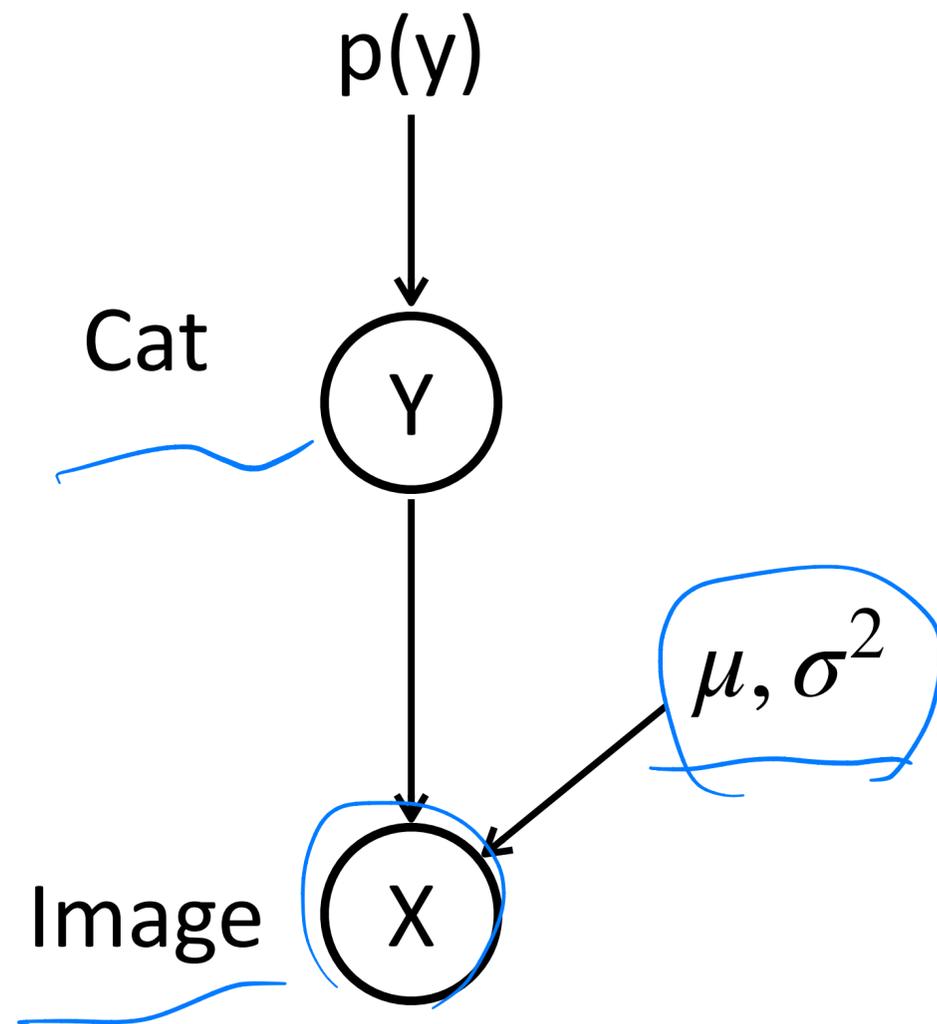


$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta | D) \\ &= \arg \max_{\theta} P(D | \theta)P(\theta)\end{aligned}$$

Bayesian statistics: there is no “parameters” in the world, all are posterior distributions to estimate

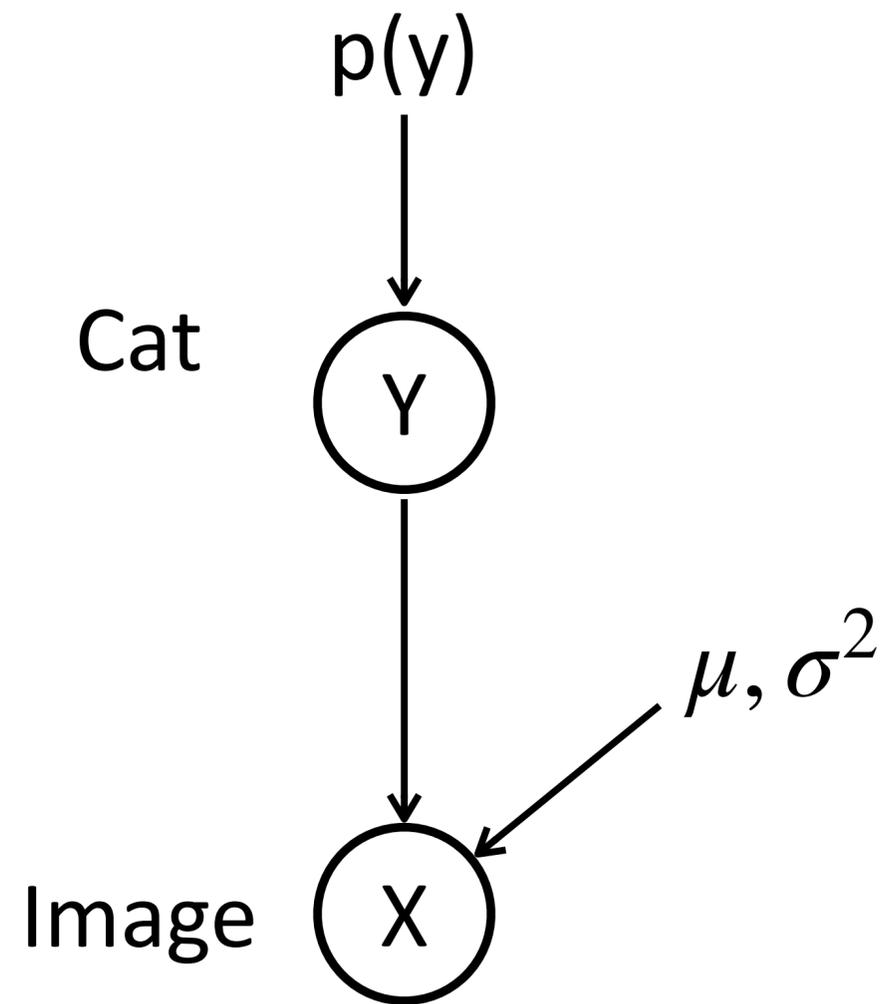
Max A Posterior (MAP) Estimation

Max A Posterior (MAP) Estimation

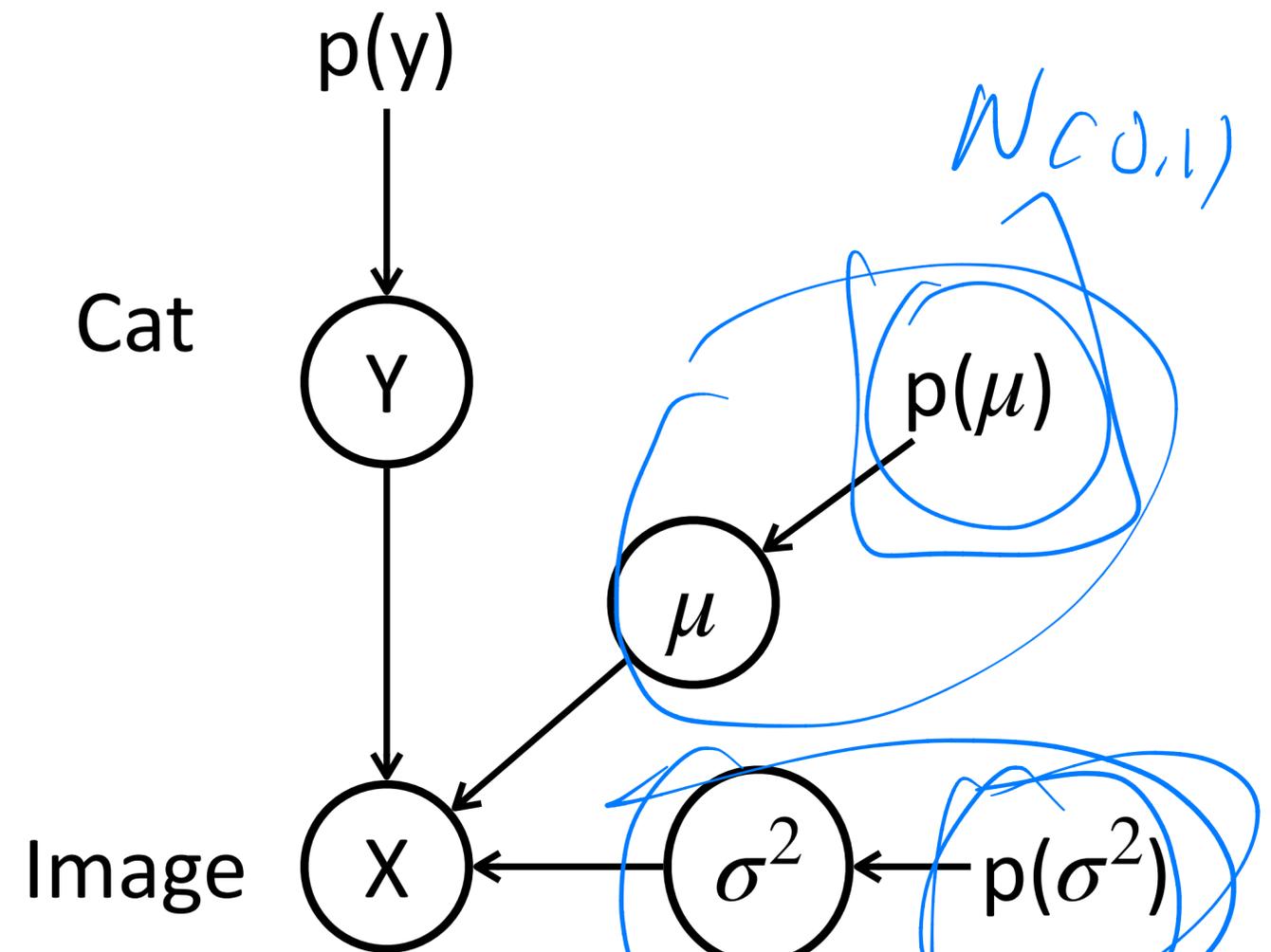


Frequentist

Max A Posterior (MAP) Estimation



Frequentist



Bayesian

How to Choose Prior

- Inject prior human knowledge to regularize the estimate
- Could learn better if data is limited

[coin]

How to Choose Prior

- Inject prior human knowledge to regularize the estimate
 - Could learn better if data is limited

● Posterior easy to compute

- Conjugate prior

exponential $P(\theta | D)$

Conjugate Prior

Conjugate Prior

If $P(\theta)$ is conjugate prior for $P(D|\theta)$, then Posterior has same form as prior

Posterior = Likelihood x Prior

$$P(\theta|D) = P(D|\theta) \times P(\theta)$$

Gauss

Gaussian

Gaussian

Conjugate Prior

If $P(\theta)$ is conjugate prior for $P(D|\theta)$, then Posterior has same form as prior

$$\text{Posterior} = \text{Likelihood} \times \text{Prior}$$
$$P(\theta|D) = P(D|\theta) \times P(\theta)$$

topic model

$P(\theta)$	$P(D \theta)$	$P(\theta D)$
Gaussian	Gaussian	Gaussian
Beta	Bernoulli	Beta
Dirichlet	Multinomial	Dirichlet

MLE vs. MAP

Maximum Likelihood estimation (MLE)

Choose value that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D|\theta)$$

Maximum *a posteriori* (MAP) estimation

Choose value that is most probable given observed data and prior belief

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max_{\theta} P(\theta|D) \\ &= \arg \max_{\theta} P(D|\theta)P(\theta)\end{aligned}$$

When are they the same?

Thank You!
Q & A