

# Machine Learning 1

Lecture 2 - Probability Theory - Maximum Likelihood - Maximum A Posteriori - Bayesian Prediction

*Erik Bekkers*



# Machine Learning 1

Recap basic probability + **independent random variables**

*Erik Bekkers*



# The Rules of Probability Theory

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

For random variables  $x \in X$  and  $y \in Y$ :

' {heads, tails}

	Discrete	Continuous
Additivity	$p(x \in A) = \sum_{x \in A} p(X = x)$	$p(x \in (a, b)) = \int_a^b p(x)dx$
Positivity	$p(X = x) \geq 0$	$p(X = x) \geq 0$
Normalization	$\sum_X p(X) = 1$	$\int_X p(x)dx = 1$
Sum Rule	$p(X) = \sum_Y p(X, Y)$	$p(x) = \int_Y p(x, y)dy$
Product Rule	$p(X, Y) = p(X   Y)p(Y)$	$p(x, y) = p(x   y)p(y)$

# Bayes Theorem

- Product rule

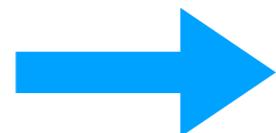
$$p(x, y) = p(x | y)p(y)$$

- Symmetry property

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

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

- Bayes rule



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

- Denominator:  $\sum_{y \in Y} p(y | x) = 1$



# Independent Random Variables

Two random variables  $X$  and  $Y$  are *independent* iff measuring  $X$  gives no information on  $Y$ , and vice versa.

- Formally:  $X$  and  $Y$  are called independent if

$$\underline{p(x,y)} = \underline{p(x)} p(y) \quad \text{for all } x \in \mathcal{X}, y \in \mathcal{Y}$$

- Equivalent to

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

$$\rightarrow p(x,y) = \underline{p(x|y)} p(y)$$

# Machine Learning 1

Lecture 2.1 - Expectation - Variance

Erik Bekkers

(Bishop 1.2.2)



# Expectations

- random variable  $x \in X$  and function  $f: X \rightarrow \mathbb{R}$

$$\mathbb{E}[f] = \mathbb{E}_{x \sim p(X)}[f(x)] = \begin{cases} \sum_{x \in X} f(x) p(x) & \text{discrete} \\ \int_X f(x) p(x) dx & \text{continuous} \end{cases}$$

"being explicit in what the expectation is taken over"

- For  $N$  points drawn from  $p(X)$ :

$$\mathbb{E}[f] \approx \frac{1}{N} \sum_{n=1}^N f(x_n)$$

- Conditional expectation:

$$\mathbb{E}[f|y] = \mathbb{E}_{x \sim p(X|Y=y)}[f(x)] = \begin{cases} \sum_{x \in X} f(x) p(x|y) \\ \int_X f(x) p(x|y) dx \end{cases}$$

# Variance

Linearity of  $\mathbb{E}$

- $\mathbb{E}[f(x) + g(x)] = \mathbb{E}[f(x)] + \mathbb{E}[g(x)]$
- $\mathbb{E}[c f(x)] = c \mathbb{E}[f(x)]$

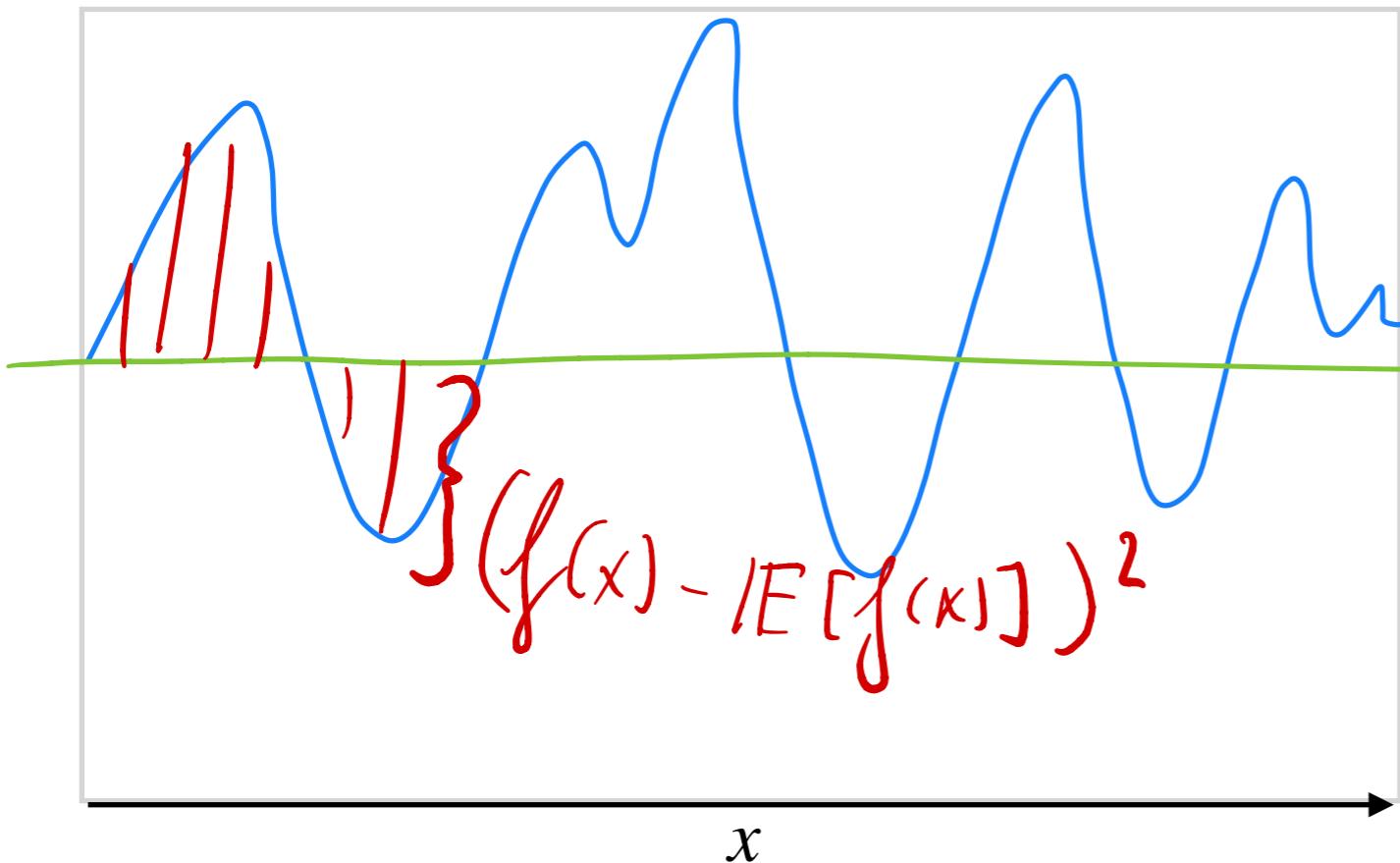
- $\mathbb{E}[c] = c$

- The expected quadratic distance between  $f$  and its mean  $\mathbb{E}[f]$

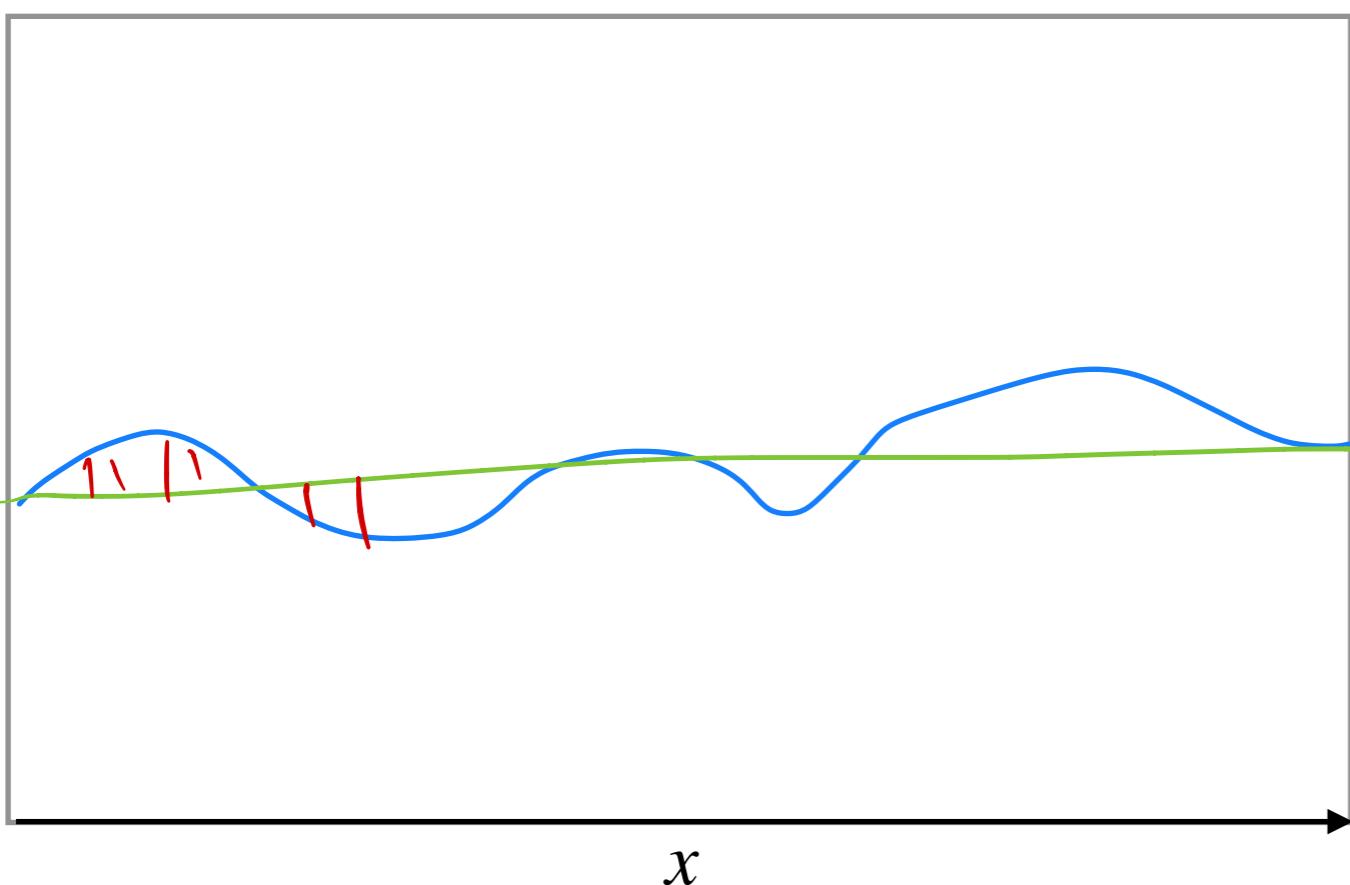
$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

$$\begin{aligned} &= \mathbb{E}[f(x)^2 - 2f(x)\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^2] \\ &= \mathbb{E}[f(x)^2] - \underbrace{\mathbb{E}[2f(x)\mathbb{E}[f(x)]]}_{\text{brace}} + \mathbb{E}[\mathbb{E}[f(x)]^2] \end{aligned}$$

$$= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$



$E[f(x)]$   
high variance



low variance

# Covariance between 2 random variables

- Measures the extent to which  $X$  and  $Y$  vary together

$$\text{cov}[x, y] = \mathbb{E}_{x,y \sim p(x,y)}[(x - \mathbb{E}[x])(y - \mathbb{E}[y])]$$

$$= \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

- Vectors of random variables  $\mathbf{x}$  and  $\mathbf{y}$ , covariance matrix:

$$\text{cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p(\mathbf{x}, \mathbf{y})}[(\underbrace{\mathbf{x} - \mathbb{E}[\mathbf{x}]}_{\mathbb{R}^{D \times 1}})(\underbrace{\mathbf{y} - \mathbb{E}[\mathbf{y}]}_{\mathbb{R}^{D \times 1}})^T] \in \mathbb{R}^{D \times D}$$

$$= \mathbb{E}[xy^T] - \mathbb{E}[x]\mathbb{E}[y]^T$$

- Define  $\text{cov}[x] := \text{cov}[x, x]$

# Covariance between 2 random variables

- Covariance between independent variables

$$\text{cov}[x, y] = \underbrace{\mathbb{E}[xy]} - \mathbb{E}[x]\mathbb{E}[y] = 0$$

$$\int \int_{X \times Y} xy p(x,y) dx dy = \int_X \left( \int_Y xy p(x) dx \right) p(y) dy = \int_X x p(x) dx \int_Y y p(y) dy = \mathbb{E}[x]\mathbb{E}[y]$$

- Note:  $\text{cov}[x, y] = 0$  does not imply  $x, y$  independent!

$$x \sim U[-1, 1] \quad p(x) = \frac{1}{2}$$

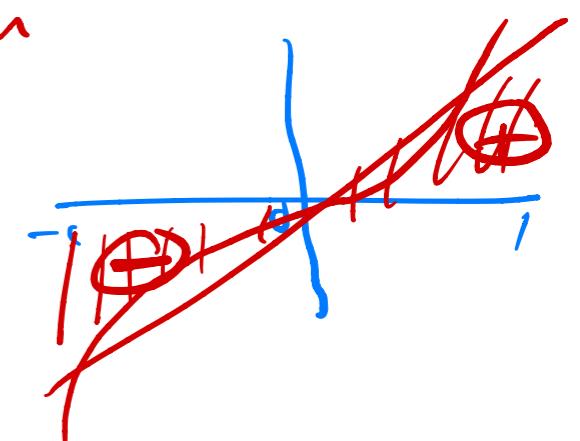
$$y = x^2$$

$$= \int x^3 p(x) dx - \int x^3 \frac{1}{2} dx = 0$$

$$\text{cov}[x, y] = \underbrace{\mathbb{E}[xy]}_{\substack{x^2 \\ 0}} - \underbrace{\mathbb{E}[x]\mathbb{E}[y]}_{\substack{0 \\ 0}} = 0$$

$x$  also odd function

because integrating odd function over symmetric domain



# Machine Learning 1

Lecture 2.2 - Gaussian Distribution

Erik Bekkers

(Bishop 1.2.4)



# Gaussian Distribution

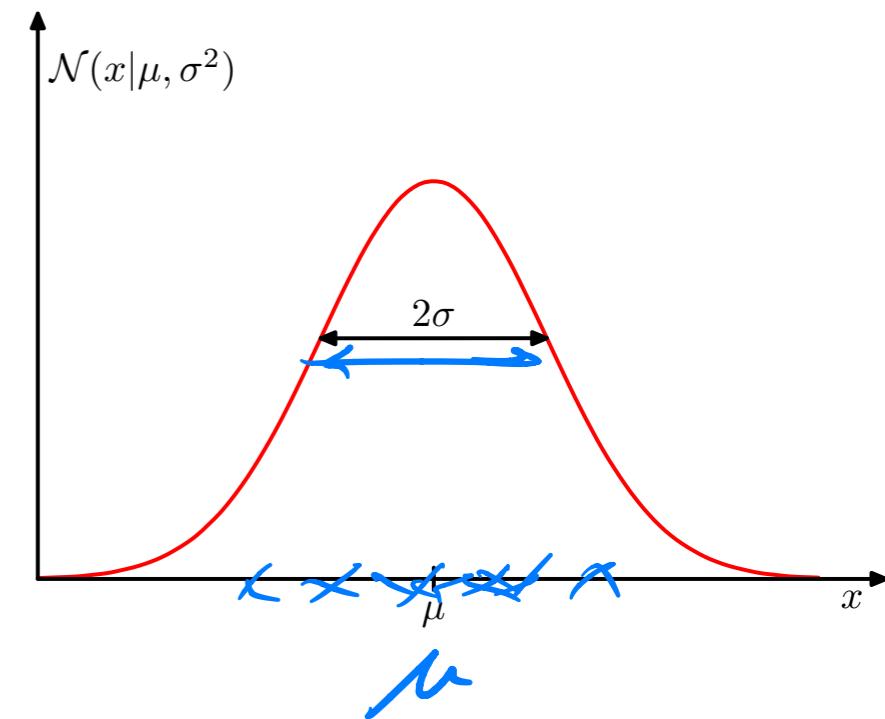
- ▶ The Normal/Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

ER  
mean

Var :  $\sigma^2$

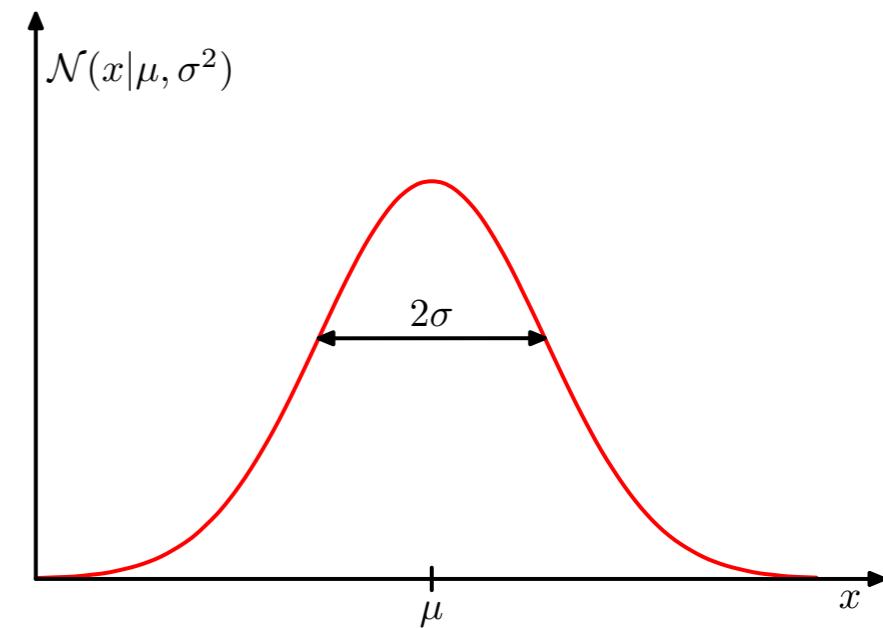
standard  
deviation :  $\sigma$



# Gaussian Distribution

- ▶ The Normal/Gaussian distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



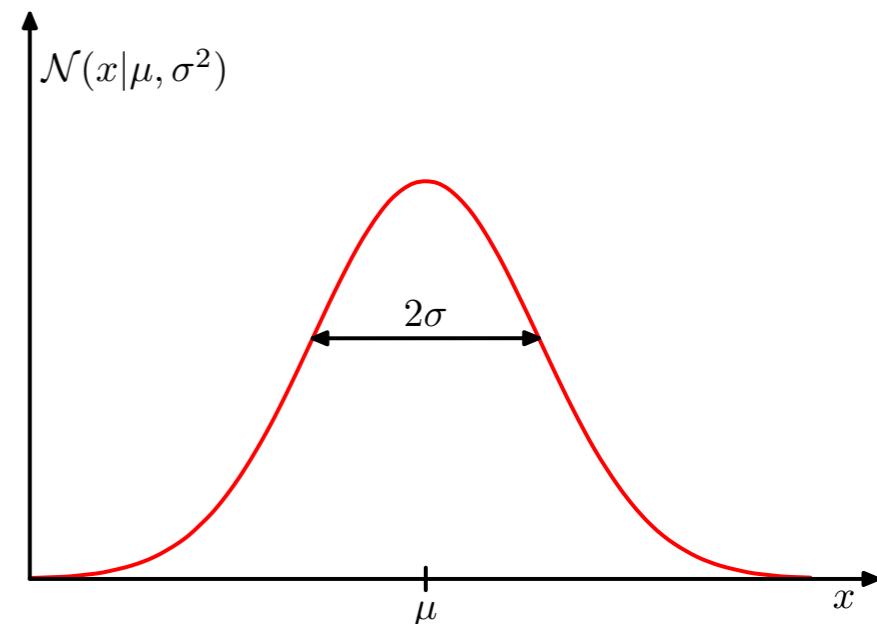
- ▶ **Mean** (of  $x \sim \mathcal{N}(x | \mu, \sigma^2)$ ):

$$\mathbb{E}_{x \sim \mathcal{N}(x | \mu, \sigma^2)}[x] = \int_{-\infty}^{\infty} x \mathcal{N}(x | \mu, \sigma^2) dx = \mu$$

# Gaussian Distribution

- The Normal/Gaussian distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



- Mean** (of  $x \sim \mathcal{N}(x | \mu, \sigma^2)$ ):

$$\begin{aligned} \mathbb{E}_{x \sim \mathcal{N}(x | \mu, \sigma^2)}[x] &= \int_{-\infty}^{\infty} x \mathcal{N}(x | \mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \end{aligned}$$

*simpl. if y to  $e^{-y^2}$*

*Step 1: Change of variables*

$$\begin{aligned} y &= \frac{1}{\sqrt{2\sigma^2}}(x - \mu) \rightarrow x = \sqrt{2\sigma^2} y + \mu \\ \frac{dy}{dx} &= \frac{1}{\sqrt{2\sigma^2}} \rightarrow dx = \sqrt{2\sigma^2} dy \end{aligned}$$

*2: Integration of odd funcs*

$$\int_{-\infty}^{\infty} y e^{-y^2} dy = 0$$

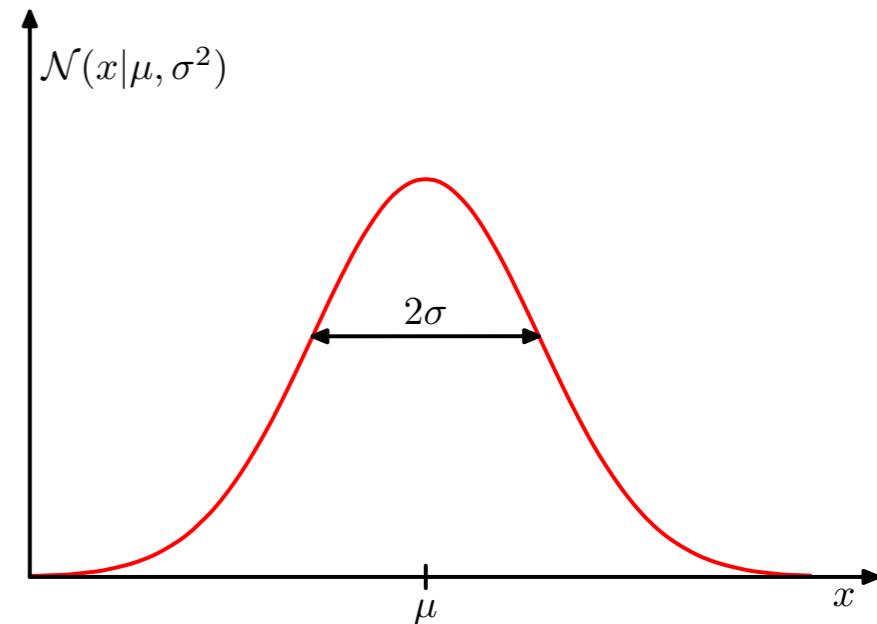
*3: Useful property:*

$$\int_{-\infty}^{\infty} e^{-y^2} dy = \sqrt{\pi}$$

# Gaussian Distribution

- ▶ The Normal/Gaussian distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



- ▶ **Mean** (of  $x \sim \mathcal{N}(x | \mu, \sigma^2)$ ):

$$\begin{aligned} \mathbb{E}_{x \sim \mathcal{N}(x | \mu, \sigma^2)}[x] &= \int_{-\infty}^{\infty} x \mathcal{N}(x | \mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\ &\stackrel{(1)}{=} \int_{-\infty}^{\infty} (\sqrt{2\sigma^2} y + \mu) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-y^2} \sqrt{2\sigma^2} dy \end{aligned}$$

*Step 1: Change of variables*

$$\begin{aligned} y &= \frac{1}{\sqrt{2\sigma^2}}(x - \mu) \rightarrow x = \sqrt{2\sigma^2} y + \mu \\ \frac{dy}{dx} &= \frac{1}{\sqrt{2\sigma^2}} \rightarrow dx = \sqrt{2\sigma^2} dy \end{aligned}$$

*2: Integration of odd funcs*

$$\int_{-\infty}^{\infty} y e^{-y^2} dy = 0$$

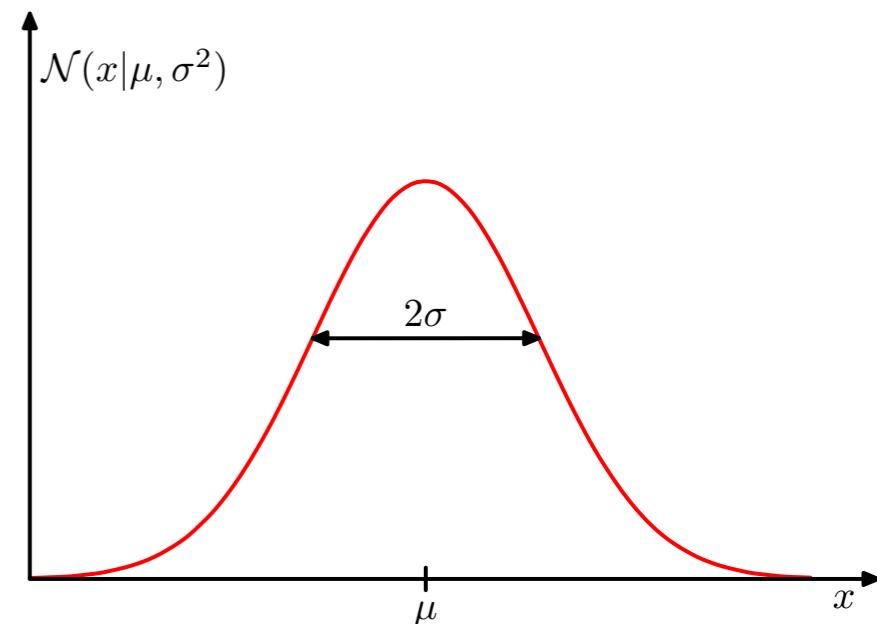
*3: Useful property:*

$$\int_{-\infty}^{\infty} e^{-y^2} dy = \sqrt{\pi}$$

# Gaussian Distribution

- The Normal/Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



- Mean** (of  $x \sim \mathcal{N}(x|\mu, \sigma^2)$ ):

$$\begin{aligned} \mathbb{E}_{x \sim \mathcal{N}(x|\mu, \sigma^2)}[x] &= \int_{-\infty}^{\infty} x \mathcal{N}(x|\mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\ &\stackrel{(1)}{=} \int_{-\infty}^{\infty} (\sqrt{2\sigma^2}y + \mu) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-y^2} \sqrt{2\sigma^2} dy \end{aligned}$$

$$= \int_{-\infty}^{\infty} \left( \cancel{\sqrt{2\sigma^2}y} + \mu \right) \frac{1}{\sqrt{\pi}} e^{-y^2} dy$$

be odd function

*Step 1: Change of variables*

$$\begin{aligned} y &= \frac{1}{\sqrt{2\sigma^2}}(x - \mu) \rightarrow x = \sqrt{2\sigma^2}y + \mu \\ \frac{dy}{dx} &= \frac{1}{\sqrt{2\sigma^2}} \rightarrow dx = \sqrt{2\sigma^2} dy \end{aligned}$$

*2: Integration of odd funcs*

$$\int_{-\infty}^{\infty} y e^{-y^2} dy = 0$$

*3: Useful property:*

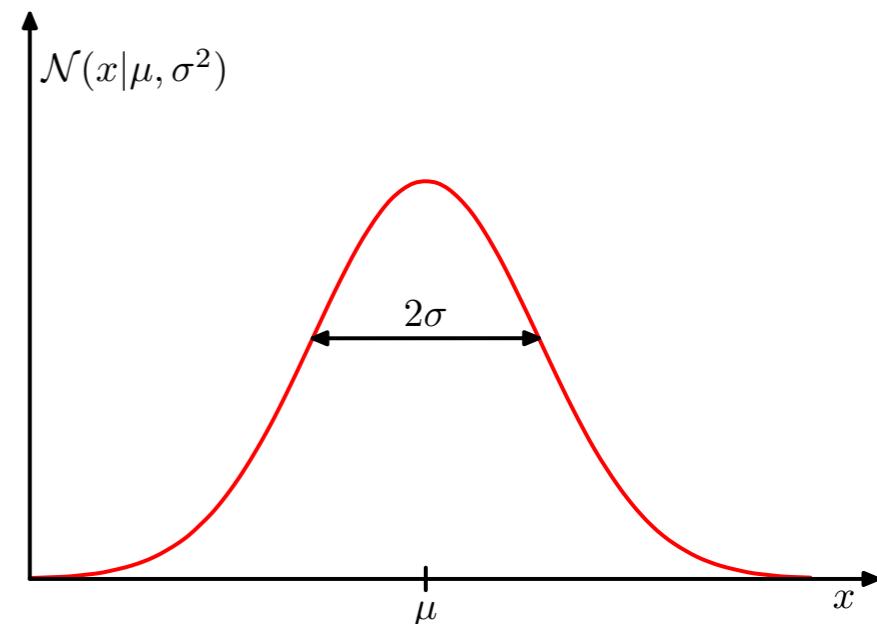
$$\int_{-\infty}^{\infty} e^{-y^2} dy = \sqrt{\pi}$$

can be derived from  
normalization property of distributions

# Gaussian Distribution

- The Normal/Gaussian distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



- Mean** (of  $x \sim \mathcal{N}(x | \mu, \sigma^2)$ ):

$$\begin{aligned}
 \mathbb{E}_{x \sim \mathcal{N}(x | \mu, \sigma^2)}[x] &= \int_{-\infty}^{\infty} x \mathcal{N}(x | \mu, \sigma^2) dx \\
 &= \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \\
 &\stackrel{(1)}{=} \int_{-\infty}^{\infty} (\sqrt{2\sigma^2} y + \mu) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-y^2} \sqrt{2\sigma^2} dy \\
 &= \int_{-\infty}^{\infty} (\sqrt{2\sigma^2} y + \mu) \frac{1}{\sqrt{\pi}} e^{-y^2} dy \stackrel{(2+3)}{=} \mu
 \end{aligned}$$

*Step 1: Change of variables*

$$\begin{aligned}
 y &= \frac{1}{\sqrt{2\sigma^2}}(x - \mu) \rightarrow x = \sqrt{2\sigma^2} y + \mu \\
 \frac{dy}{dx} &= \frac{1}{\sqrt{2\sigma^2}} \rightarrow dx = \sqrt{2\sigma^2} dy
 \end{aligned}$$

*2: Integration of odd funcs*

$$\int_{-\infty}^{\infty} y e^{-y^2} dy = 0$$

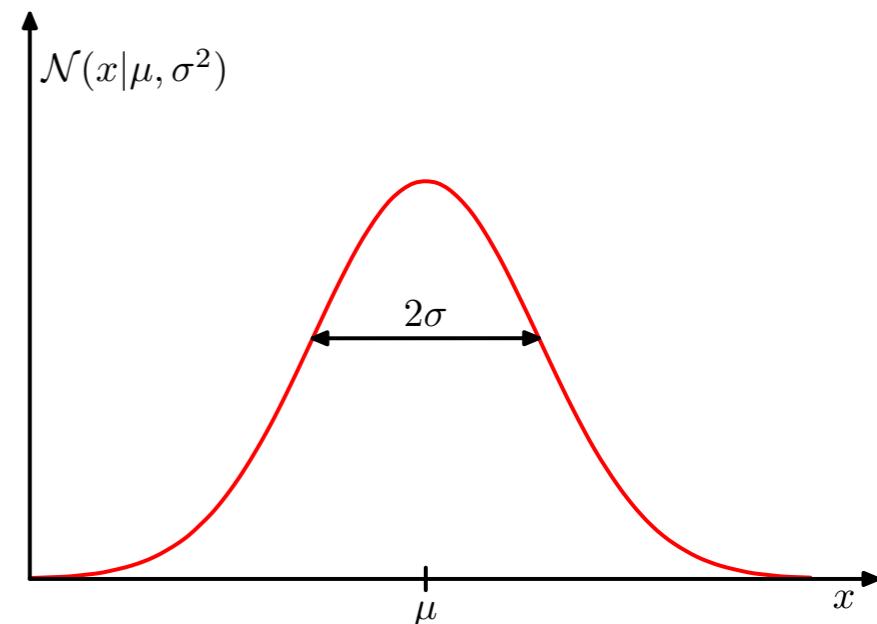
*3: Useful property:*

$$\int_{-\infty}^{\infty} e^{-y^2} dy = \sqrt{\pi}$$

# Gaussian Distribution

- ▶ The Normal/Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



- ▶ **Variance** (of  $x \sim \mathcal{N}(x|\mu, \sigma^2)$ ):

$$\begin{aligned} \text{var}[x] &= \int_{-\infty}^{\infty} (x - \mu)^2 \mathcal{N}(x|\mu, \sigma^2) dx \\ &= \int_{-\infty}^{\infty} (x - \mu)^2 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} dx \end{aligned}$$

... (using 1 + 2 + 4)

=  $\sigma^2$

see Ch 2.3

*Step 1: Change of variables*

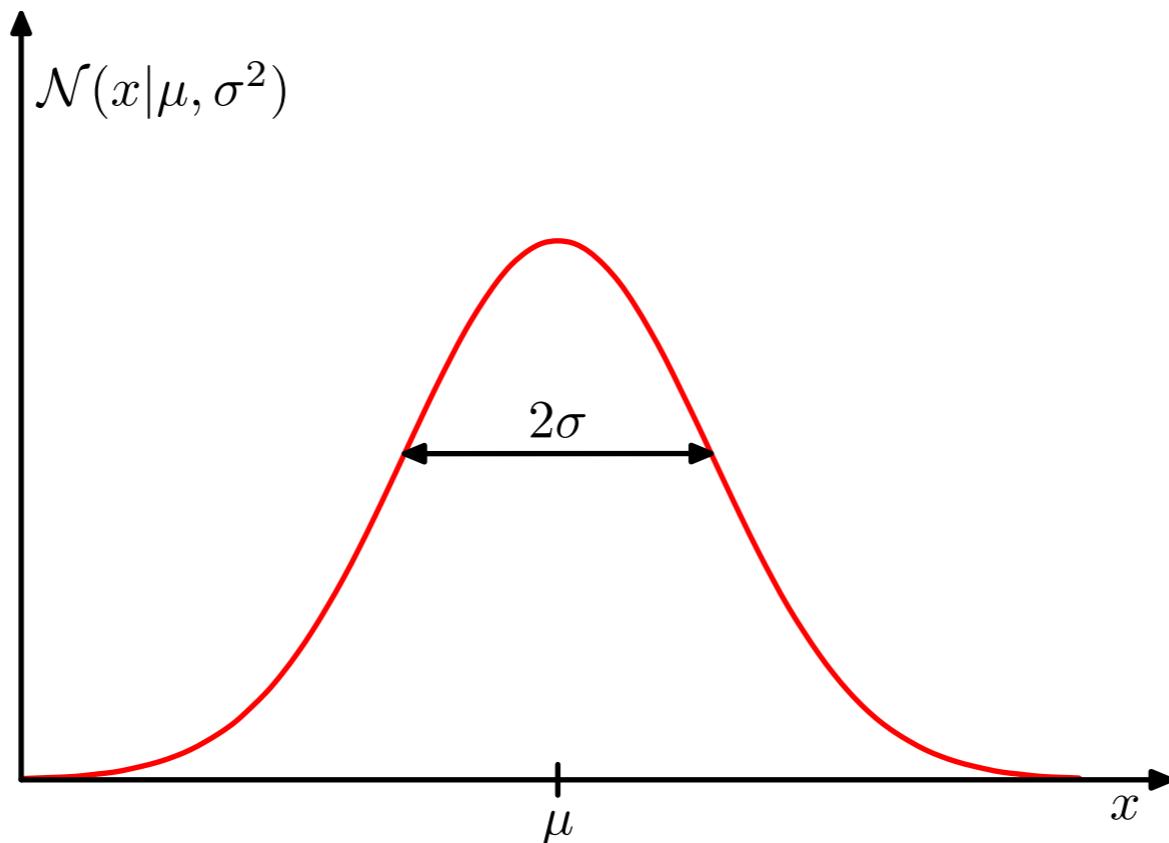
$$\begin{aligned} y &= \frac{1}{\sqrt{2\sigma^2}}(x - \mu) \rightarrow x = \sqrt{2\sigma^2}y + \mu \\ \frac{dy}{dx} &= \frac{1}{\sqrt{2\sigma^2}} \rightarrow dx = \sqrt{2\sigma^2} dy \end{aligned}$$

*4: Another convenient trick:*

$$\begin{aligned} \int_{-\infty}^{\infty} x^2 e^{-ax^2} dx &= -\frac{\partial}{\partial a} \int_{-\infty}^{\infty} e^{-ax^2} dx \\ &= \frac{\partial}{\partial a} \sqrt{\frac{\pi}{a}} = \frac{1}{2} \sqrt{\frac{\pi}{a^3}} \end{aligned}$$

# Gaussian Distribution

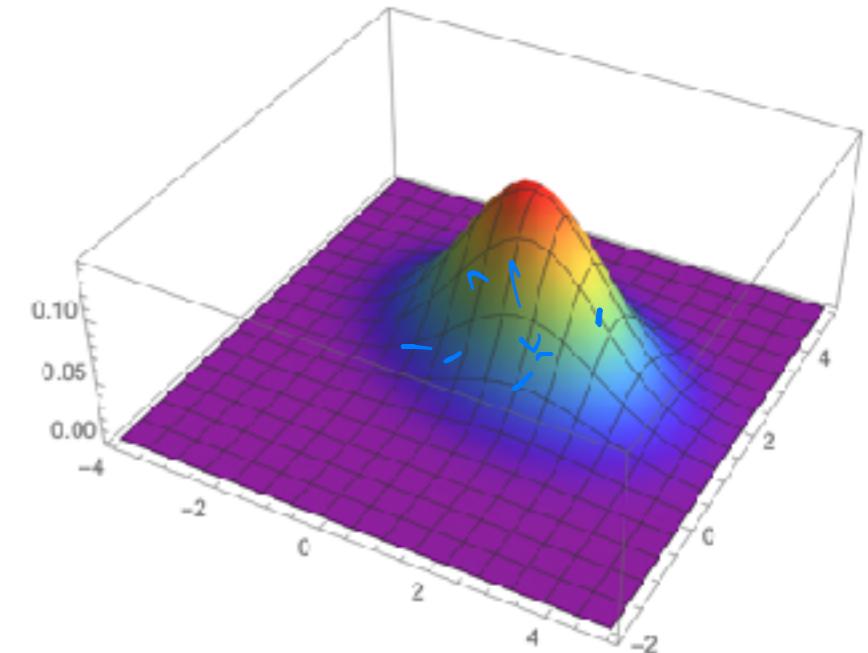
$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



$$x \sim \mathcal{N}(x | \mu, \sigma^2) : \left\{ \begin{array}{l} \mathbb{E}[x] = \mu \\ \text{Var}[x] = \sigma^2 \end{array} \right.$$

# Multivariate Gaussian Distribution

- $D$ -dimensional vector  $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$
- $\mathcal{N}(\mathbf{x} | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^D |\Sigma|}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1} (\mathbf{x}-\mu)}$   $|\Sigma| = \det \Sigma$
- $\Sigma = \text{Cov}[\underline{\mathbf{x}}, \underline{\mathbf{x}}] \in \mathbb{R}^{D \times D}$
- $\mathbb{E}[\mathbf{x}] = \underline{\mu}$



1. Substitution  $\mathbf{y} = (\mathbf{x} - \mu)$

2. Normalization factor:  $\int_{\mathbb{R}^D} e^{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x}} = \frac{2\pi^{D/2}}{|\mathbf{A}|^{1/2}}$  (since  $\int_{\mathbb{R}^D} \mathcal{N}(\mathbf{x} | \mu, \Sigma) = 1$ )

$N(x | \mu, \Sigma)$



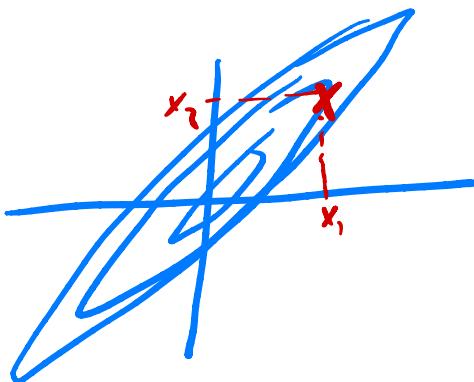
isotropic

$$\Sigma = I \approx \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix}$$

isotropic

$x_1$  and  $x_2$   
are not correlated

extremely  
anisotropic



$x_1$  and  $x_2$   
are highly  
correlated

knowing  $x_1$  is  
knowing  $x_2$

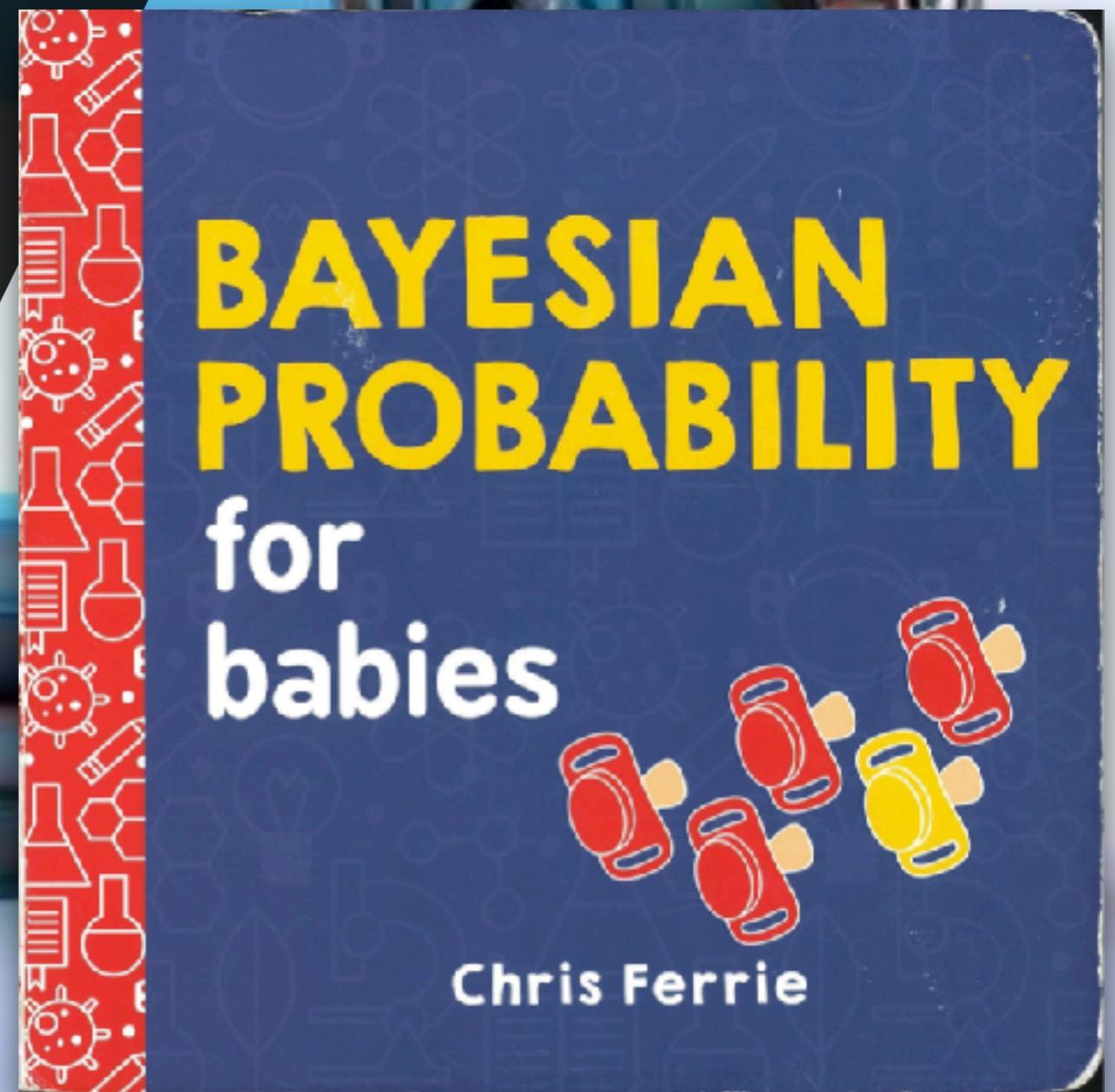
# Machine Learning 1

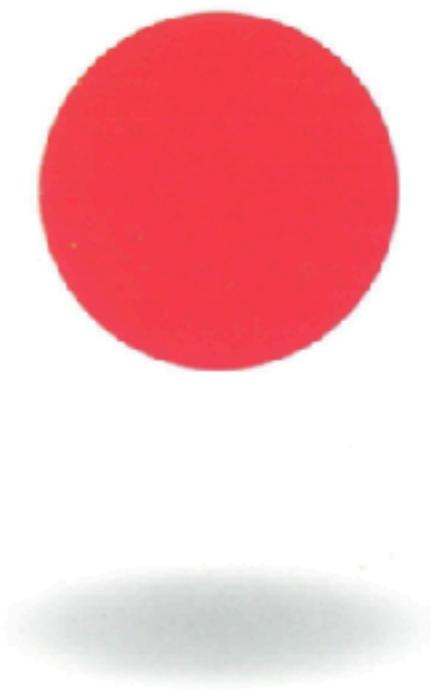
## Example Bayes Rule

Erik Bekkers

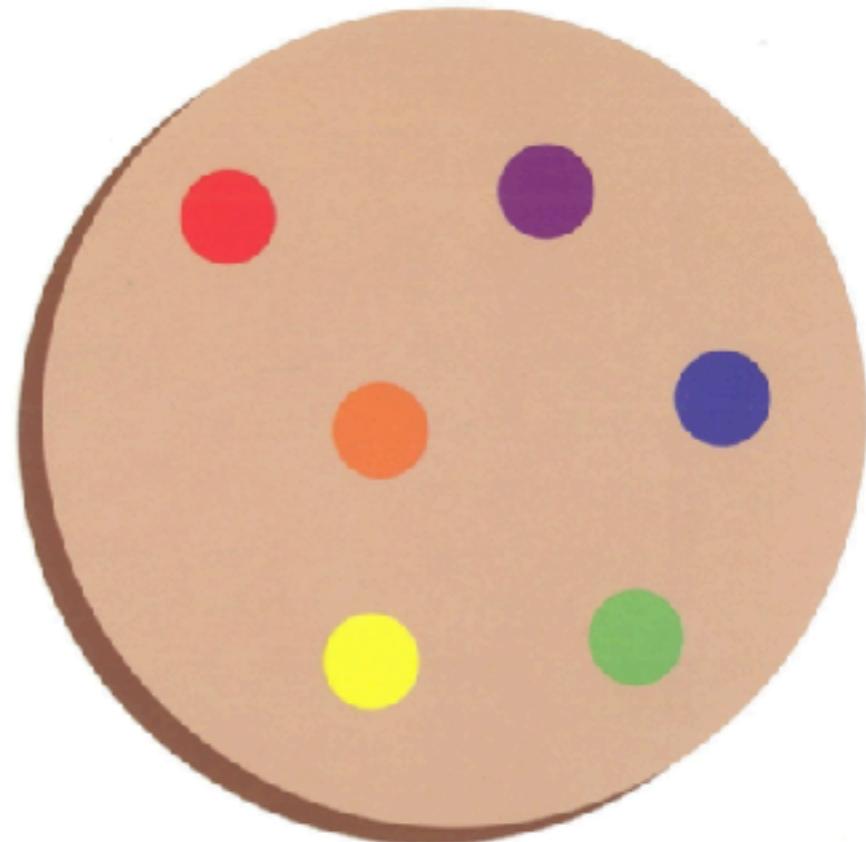
(Bishop 1.2.0 - 1.2.1)

Getting into model  
estimation





This is a ball.



The ball is a piece of candy on a cookie. Yum!



$C=1$

Some cookies have candy.



$C=0$

Some don't.

Cookie

$C = \{ \text{candy}, \text{no-candy} \}$

{ | , 0 }



$B=0$

Take a bite. It has no candy.

$$B = \{1, 0\}$$

= {with candy, no candy}



$B=0$

Did it come from a candy cookie?



Either it came from a candy cookie,

likelihood empty bag came  
from candy cookie

$$p(B=0 | C \geq 1)$$

posterior prob

$$p(C \geq 1 | B=0)$$

or it didn't. What are the chances?

$$p(B=0 | C=0)$$

$$p(C=0 | B=0)$$



If the cookie had no candy,  
then every bite would have no candy.

$$\Pr(\text{bite} | \text{no candy}) = 1$$

The probability of a no-candy bite,  
given a no-candy cookie, is 1.

$$p(B=0 | C=0) = 1$$

"likelihood under the no candy  
model"



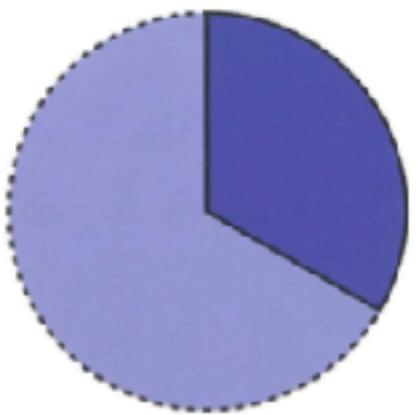
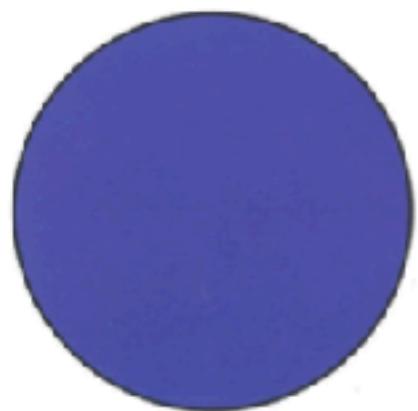
If the cookie had candy, then very few bites would have no candy.

$$\Pr(\text{bite} \mid \text{candy}) = \frac{1}{3}$$

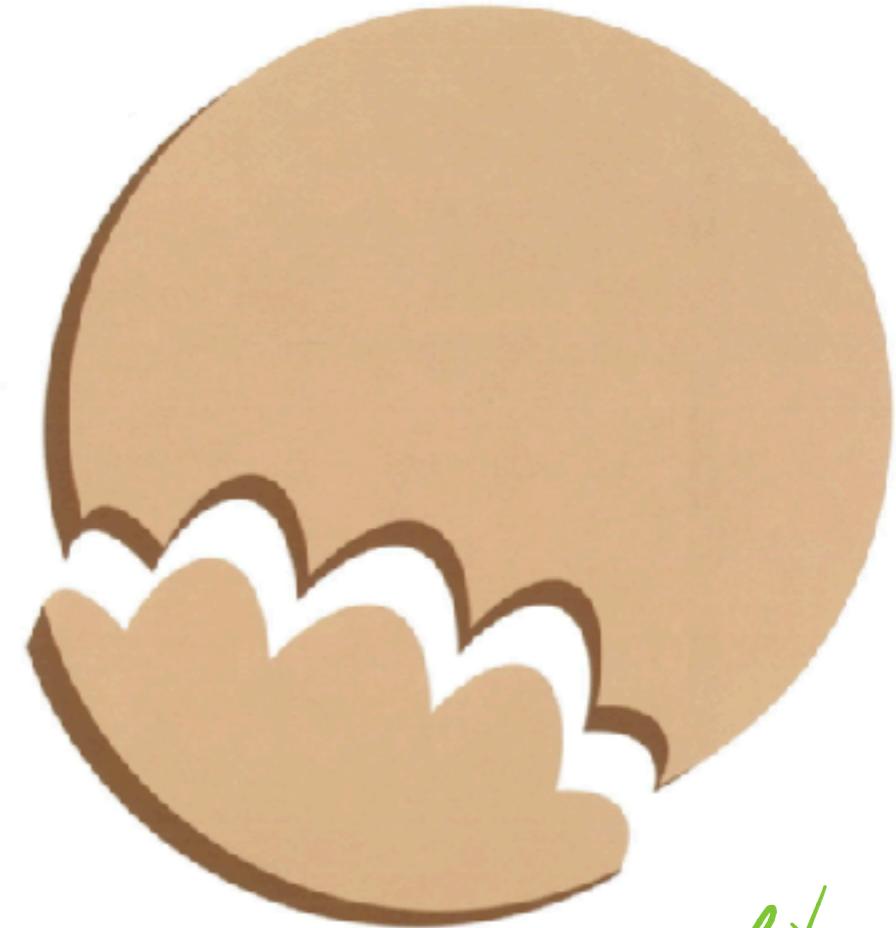
The probability of a no-candy bite, given a candy cookie, is  $1/3$ .

$$\Pr(\text{Bite} \mid \text{Candy}) = \frac{1}{3}$$

$$\Pr(\text{bite} | \text{no candy}) > \Pr(\text{bite} | \text{candy})$$

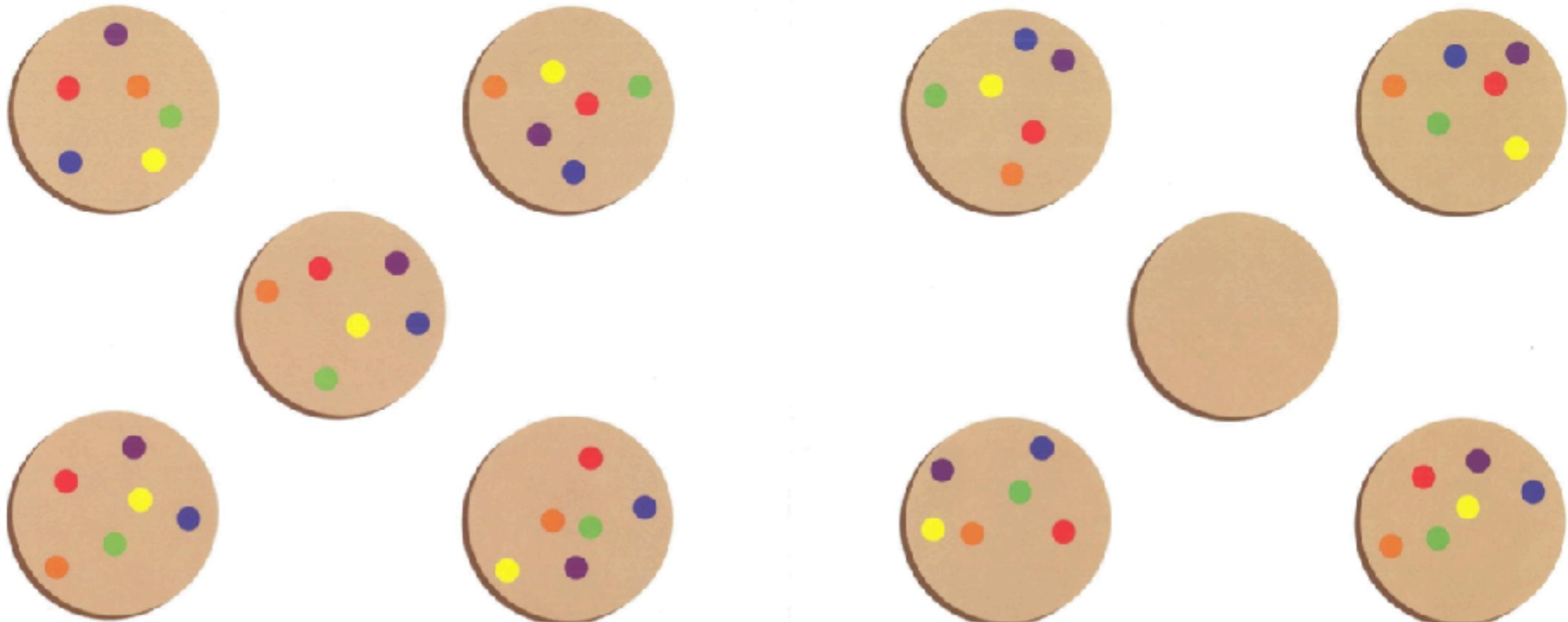


1 is greater than 1/3.



So the no-candy bite probably came from a no-candy cookie!

*likely*



But what if we knew there were 10 cookies,

and all had candy but one?

prior  $p(C=1) = \frac{9}{10}$

$$p(C=0) = \frac{1}{10}$$



Take a bite of each.

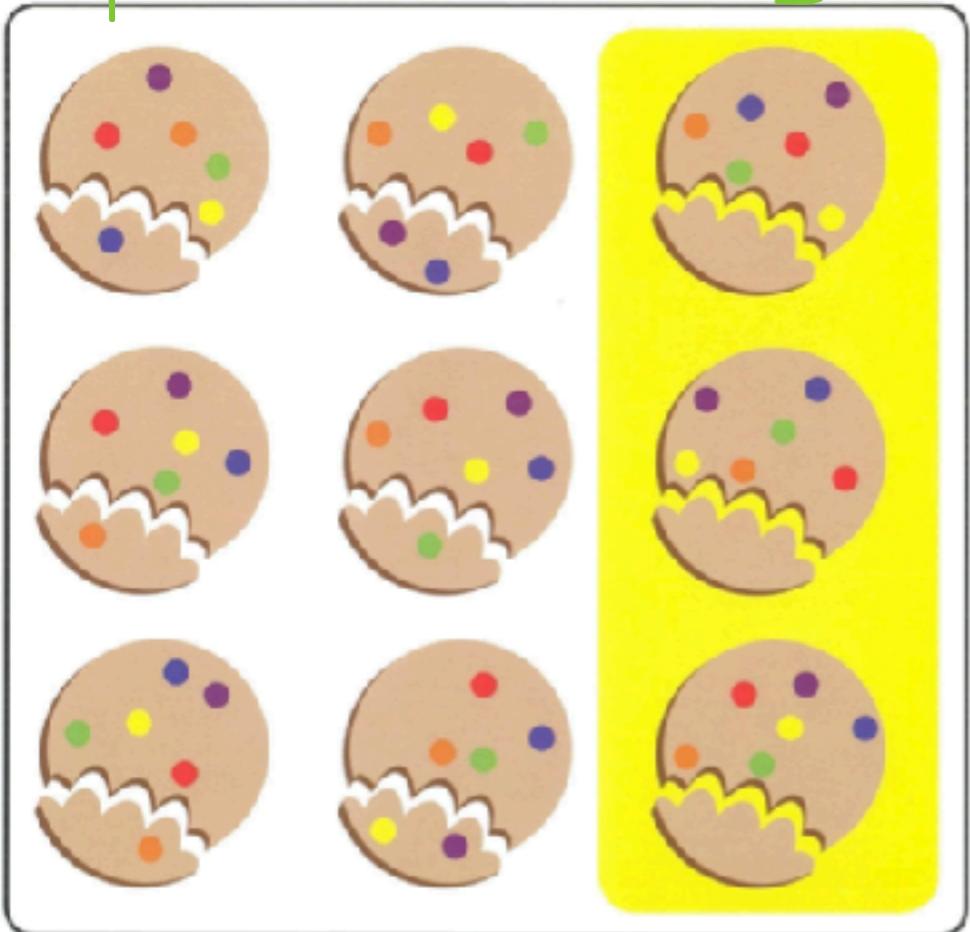


There are 4 no-candy bites.  
3 bites are from candy cookies.  
1 bite is from a no-candy cookie.

$$p(B=0) = \frac{4}{10}$$

$$p(B=1) = \frac{6}{10}$$

$$p(B=0 | C=1) = \frac{1}{3}$$



$$\Pr(\text{Cookie with no candy} | \text{Bite taken}) = \frac{3}{4}$$

1/3 of the candy cookies have a no-candy bite.

The probability of a candy cookie with a no-candy bite is 3/4.

Posterior prob

$$p(CC=1 | B=0) = \frac{p(B=0 | C=1) \cdot p(C=1)}{p(B=0)}$$

$$= \frac{\frac{1}{3} \cdot \frac{9}{10}}{\frac{4}{10}} = \frac{3}{4}$$



This is the **prior distribution** of cookies.

$$p(C=1) = \frac{9}{10}$$



This is the **posterior distribution** of cookies.

$$p(C=1 | B=0) = \frac{3}{4}$$



This bite probably came from a candy cookie!

Classification: pick the most probable option



# Machine Learning 1

Lecture 2.3 - Maximum Likelihood

Erik Bekkers

(Bishop 1.2.3 - 1.2.5)



***Three Statistical Learning Principles:***

- Maximum Likelihood**
- Maximum A Posteriori**
- Bayesian Prediction**





# Maximum Likelihood Estimation

- **Objective:** Find maximum likelihood solution  $\mathbf{w}_{ML}$

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} p(D | \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_i^N p(x_i | \mathbf{w})$$

*L1*

- **Problem:** numerical underflow/overflow

**Solution:**

# Maximum Likelihood Estimation

- **Objective:** Find maximum likelihood solution  $\mathbf{w}_{ML}$

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} p(D | \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_i^N p(x_i | \mathbf{w})$$

- **Problem:** numerical underflow/overflow

$$\log(a \cdot b) = \log(a) + \log(b)$$

- **Solution:** Maximize log-likelihood instead:

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_i^N p(x_i | \mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \log \prod_i^N p(x_i | \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^N \log p(x_i | \mathbf{w})$$

- Or minimize error function negative log-likelihood function:

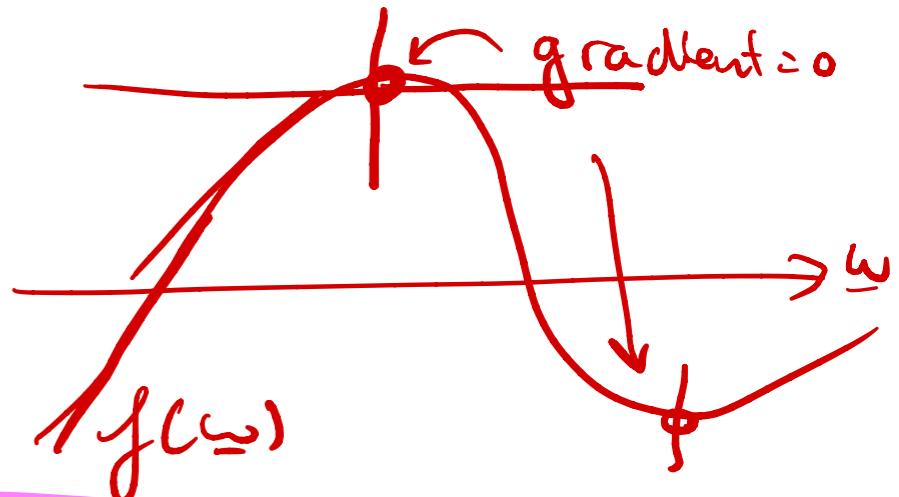
$$E(D | \mathbf{w}) = -\log p(D | \mathbf{w})$$

if we want to frame it  
as a minimization problem

# Why log-trick?

- Solution to argmax problem is obtained by solving:

$$\underset{\mathbf{w}}{\operatorname{argmax}} \ p(D | \mathbf{w}) \Leftrightarrow \text{Solve: } \frac{\partial}{\partial \mathbf{w}} p(D | \mathbf{w}) = 0$$



- We can take the log as  $p(x_i | \mathbf{w}) > 0$

$$\underset{\mathbf{w}}{\operatorname{argmax}} \ \log p(D | \mathbf{w}) \Leftrightarrow \text{Solve: } \frac{\partial}{\partial \mathbf{w}} \log p(D | \mathbf{w}) = 0$$

Same problem? Yes!

$$\frac{\partial}{\partial \mathbf{w}} \log p(D | \mathbf{w}) = \cancel{\frac{1}{p(D | \mathbf{w})}}$$

recall  $\frac{\partial}{\partial x} \log x = \frac{1}{x}$

Since  $p(D | \mathbf{w})$  will never be zero, thus  $\frac{\partial}{\partial \mathbf{w}} p(D | \mathbf{w})$  has to be.

# ML Estimator for Gaussian Distributions

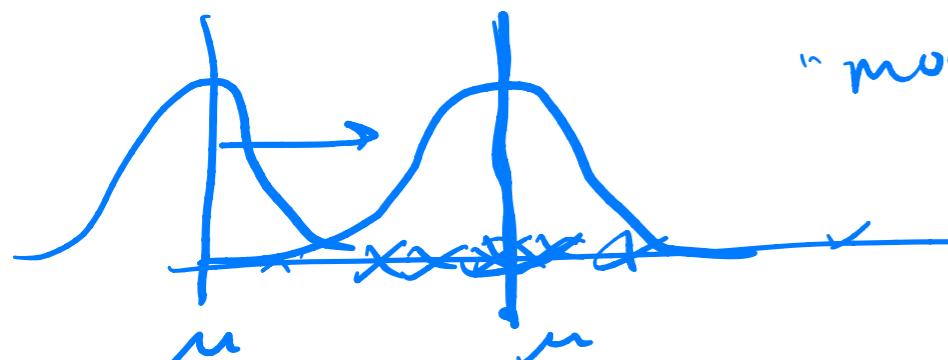
$$\log a \cdot b = \log a + \log b$$
$$\log a^n = n \log a$$

- I.i.d. Gaussian distributed real variables  $D = (x_1, x_2, \dots, x_N)$

$$p(x | w) = \mathcal{N}(x | \mu, \sigma^2) \rightarrow p(D | w) = (2\pi\sigma^2)^{-N/2} \prod_{i=1}^N \exp \left[ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

- Log likelihood

$$\log p(D | \mu, \sigma^2) = -\frac{N}{2} \log (2\pi\sigma^2) - \sum_{i=1}^N \frac{1}{2\sigma^2} (x_i - \mu)^2$$



"move/find  $\mu$  as to maximize  
(log) likelihood"

- Estimate model parameters:

$$\mu_{ML}, \sigma_{ML} = \operatorname{argmax}_{\mu, \sigma} \log p(D | \mu, \sigma^2)$$

# ML Estimator for Gaussian Distributions

- log likelihood:

$$\log p(D | \mu, \sigma^2) = -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

1. compute derivative  
2. set to zero  
3. solve equation

- Maximum Likelihood solution for  $\mu$ :

$$\begin{aligned}\frac{\partial}{\partial \mu} \log p(D | \mu, \sigma^2) &= -\frac{1}{2\sigma^2} \sum_{i=1}^N \frac{d}{d\mu} (x_i - \mu)^2 = 0 \\ \Rightarrow -\frac{1}{\sigma^2} \sum_{i=1}^N 2(x_i - \mu) &= 0 \\ \Rightarrow \sum_{i=1}^N x_i - \sum_{i=1}^N \mu &= 0 \\ \Rightarrow \sum_{i=1}^N x_i - N\mu &= 0 \Rightarrow \mu = \frac{1}{N} \sum_{i=1}^N x_i\end{aligned}$$

$\mu_{ML}$  is Sample mean:  
$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

# ML Estimator for Gaussian Distributions

- log likelihood:

$$\log p(D | \mu, \sigma^2) = -\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

- Maximum Likelihood solution for  $\sigma^2$ :

$$\frac{\partial}{\partial \sigma^2} \log p(D | \mu, \sigma^2) = -\frac{N}{2} \cancel{\frac{1}{2\pi\sigma^2}} \cdot \cancel{2\pi} + \frac{1}{2\sigma^4} \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$\Rightarrow -N\sigma^2 + \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$\Rightarrow \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

$\sigma_{ML}^2$  is sample variance:

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

# ML Estimator for Gaussian Distributions

- ▶ How well do the ML estimates represent the true parameters?

$$p(D | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{N/2}} \prod_{i=1}^N \exp \left[ -\frac{1}{2\sigma^2}(x_i - \mu)^2 \right]$$

- ▶ Drawing multiple datasets  $D \sim p(D | \mu, \sigma^2)$ , what is the expected value of  $\mu_{ML}$ ?

$$D_1 = \{x_1, \dots, x_N\}, D_2 = \{x_1, \dots, x_N\}$$

- ▶ ML estimate of the mean:

$$\mathbb{E}_{D \sim p(D | \mu, \sigma^2)}[\mu_{ML}] = \mathbb{E}_{D \sim p(D | \mu, \sigma^2)} \left[ \frac{1}{N} \sum_{i=1}^N x_i \right]$$

expectation over  
datasets

- ▶ Bias of estimator:

$$\mathbb{E}[\mu_{ML}] - \mu =$$

# ML Estimator for Gaussian Distributions

- ▶ How well do the ML estimates represent the true parameters?

$$p(D | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{N/2}} \prod_{i=1}^N \exp \left[ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

- ▶ Drawing multiple datasets  $D \sim p(D | \mu, \sigma^2)$ , what is the expected value of  $\mu_{ML}$ ?
- ▶ ML estimate of the mean:

$$\mathbb{E}_{D \sim p(D | \mu, \sigma^2)} [\mu_{ML}] = \mathbb{E}_{D \sim p(D | \mu, \sigma^2)} \left[ \frac{1}{N} \sum_{i=1}^N x_i \right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{x \sim p(x | \mu, \sigma^2)} [x_i] = \frac{1}{N} \sum_{i=1}^N \mu = \mu$$

- ▶ Bias of estimator:

$$\mathbb{E}[\mu_{ML}] - \mu = O(\cdot)$$

# ML Estimator for Gaussian Distributions

- ▶ How well do the ML estimates represent the true parameters?

$$p(D | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{N/2}} \prod_{i=1}^N \exp \left[ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

- ▶ Drawing multiple datasets  $D \sim p(D | \mu, \sigma^2)$ , what is the *expected value* of  $\mu_{ML}$ ?

- ▶ ML estimate of the mean:

$$\mathbb{E}_{D \sim p(D | \mu, \sigma^2)} [\mu_{ML}] = \mathbb{E}_{D \sim p(D | \mu, \sigma^2)} \left[ \frac{1}{N} \sum_{i=1}^N x_i \right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{x \sim p(x | \mu, \sigma^2)} [x_i] = \frac{1}{N} \sum_{i=1}^N \mu = \mu$$

- ▶ Bias of estimator ML estimator for  $\mu$ :

$$\mathbb{E}[\mu_{ML}] - \mu = \textcircled{O}$$

# ML Estimator for Gaussian Distributions

- › ML estimate of the variance:

$$\mathbb{E}_{D \sim p(D|\mu, \sigma^2)}[\sigma_{ML}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right]$$

# ML Estimator for Gaussian Distributions

- ML estimate of the variance:

$$\mathbb{E}_{D \sim p(D|\mu, \sigma^2)}[\sigma_{ML}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right]$$

# ML Estimator for Gaussian Distributions

- ML estimate of the variance:

$$\mathbb{E}_{D \sim p(D|\mu, \sigma^2)}[\sigma_{ML}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ x_i^2 - \frac{2x_i}{N} \sum_{n=1}^N x_n + \frac{1}{N^2} \sum_{m=1}^N \sum_{n=1}^N x_m x_n \right]$$

*expand*

# ML Estimator for Gaussian Distributions

- ML estimate of the variance:

$$\mathbb{E}_{D \sim p(D|\mu, \sigma^2)}[\sigma_{ML}^2] = \mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ \left( x_i - \frac{1}{N} \sum_{n=1}^N x_n \right)^2 \right] = \frac{1}{N} \sum_{i=1}^N \mathbb{E} \left[ x_i^2 - \frac{2x_i}{N} \sum_{n=1}^N x_n + \frac{1}{N^2} \sum_{m=1}^N \sum_{n=1}^N x_m x_n \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \left\{ \mathbb{E}[x_i^2] - \frac{2}{N} \sum_{n=1}^N \mathbb{E}[x_i x_n] + \frac{1}{N^2} \sum_{m=1}^N \sum_{n=1}^N \mathbb{E}[x_m x_n] \right\}$$

$$= \dots = \frac{N-1}{N} \sigma^2$$

$$\mathbb{E}[x_i x_j] = \begin{cases} \mu^2 + \sigma^2 & \text{if } i = j \\ \mu^2 & \text{if } i \neq j \end{cases} \quad \leftarrow \text{derive from known covariances!}$$

# Values for $\mathbb{E}[x_i x_j]$

- $i = j$ :

$$\text{Var}[x_i] = \mathbb{E}[x_i^2] - \mathbb{E}[x_i]^2 = \sigma^2$$

$$\Rightarrow \mathbb{E}[x_i x_j] = \mu^2 + \sigma^2$$

- $i \neq j$ :

$$\text{Cov}[x_i, x_j] = \mathbb{E}[x_i x_j] - \mathbb{E}[x_i] \mathbb{E}[x_j] = 0$$

$$\Rightarrow \mathbb{E}[x_i x_j] = \mu^2$$

# ML Estimator for Gaussian Distributions (VI)

- For data generated from

$$p(D | \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{N/2}} \prod_{i=1}^N \exp \left[ -\frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

- ML gives biased estimator for  $\sigma^2$

$$\mathbb{E}[\sigma_{ML}^2] = \frac{N-1}{N} \sigma^2$$

apparently

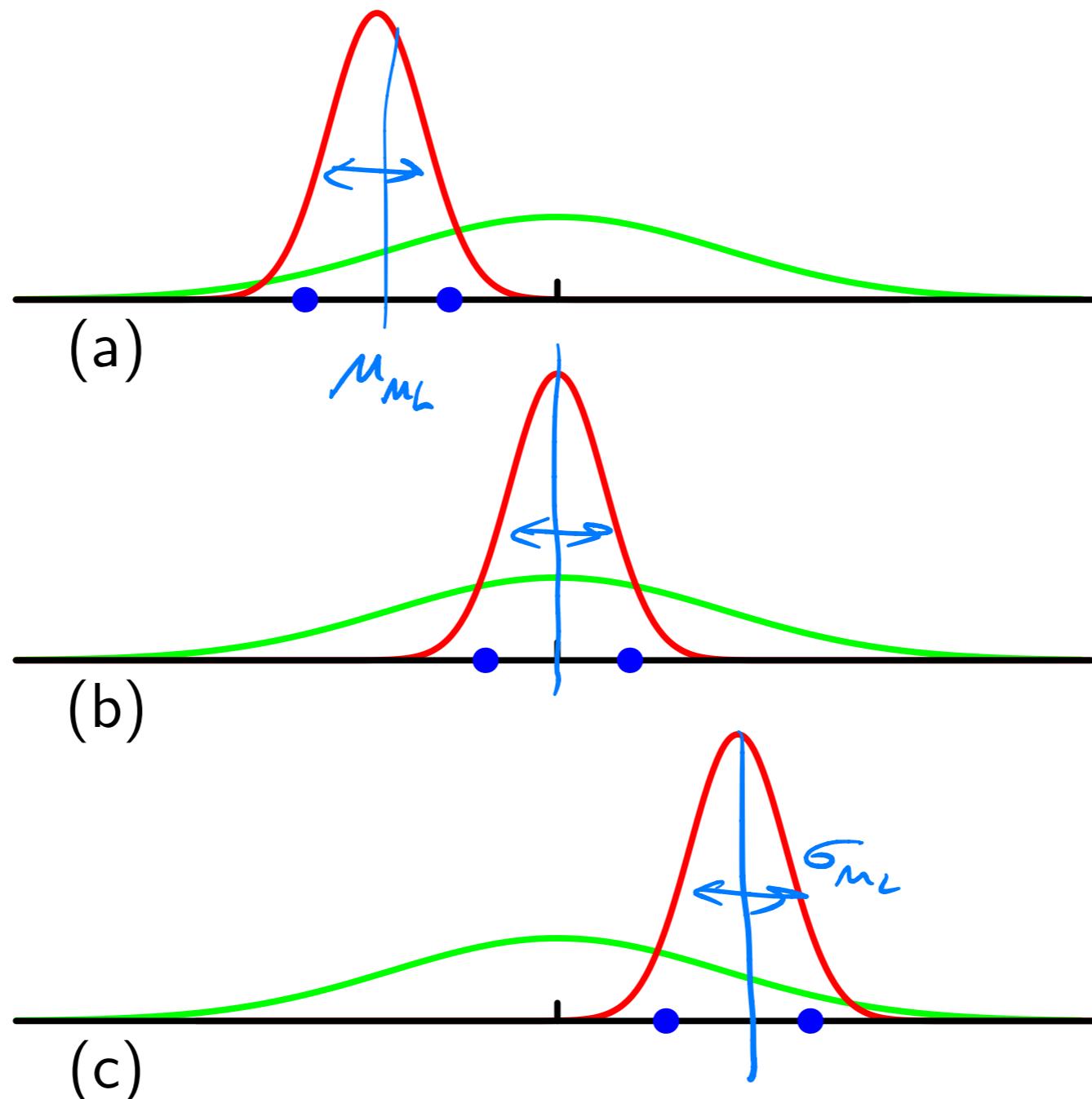
$$\mathbb{E}[(\hat{\sigma}_{ML}^2)] \neq \sigma^2$$

- Unbiased variance estimator:

$$\tilde{\sigma}^2 = \frac{N}{N-1} \sigma_{ML}^2$$

simply correct for it

# Biased Maximum Likelihood Estimator



**Figure:** Bias in ML estimator for variance (Bishop 1.15)